

## (12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-272699

(43) 公開日 平成11年(1999)10月8日

BD

識別記号  
./30  
17/27

FI  
G 0 6 F 15/401 3 2 0 A  
15/20 5 5 0 A  
15/38 D  
15/40 3 7 0 A

審査請求 未請求 請求項の数17 O L (全 74 頁)

(21) 出願番号 特願平10-72724

(22) 出願日 平成10年(1998)3月20日

(71) 出願人 000005223

富士通株式会社  
神奈川県川崎市中原区上小田中4丁目1番  
1号

(72) 発明者 仲尾 由雄

神奈川県川崎市中原区上小田中4丁目1番  
1号 富士通株式会社内

(74) 代理人 弁理士 大昔 義之 (外1名)

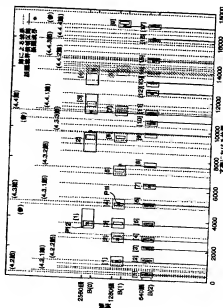
(54) 【発明の名称】 文書要約装置およびその方法

(57) 【要約】

【課題】 文章一般に見られる現象をもとに、文書中の話題構成を自動的に認定して、話題構成に対応する要約を作成することが課題である。

【解決手段】 文書要約装置は、異なる大きさの複数の窓を用いて、文書中の各位置における語彙的結束度を計算し、各話題の階層毎に話題境界の候補区間を求める。次に、異なる階層の候補区間を順に統合していくことで、各階層毎に話題境界を認定する。そして、要約作成対象の話題のまとまりと、それを含む大きな話題のまとまりとの関係に基づき、重要文を抽出して要約を作成する。

話題構成の第1の認定結果を示す図



## 【特許請求の範囲】

【請求項 1】 与えられた文書中の話題の階層的構成を認定する構成認定手段と、

認定された各話題に関する重要語を抽出する抽出手段と、

前記重要語の出現状況に応じて、各話題のまとまりから重要文を選択し、該重要文を用いて要約を生成する選択手段と、

前記要約を出力する出力手段とを備えることを特徴とする文書要約装置。

【請求項 2】 前記構成認定手段は、前記文書中の各位置の近傍領域における語彙的結束度を求め、該結束度に基づいて話題境界を認定し、該近傍領域の大きさを段階的に縮小しながら話題境界の認定を繰り返すことで、大きな話題のまとまりから小さな話題のまとまりに至る話題の階層的構成を認定することを特徴とする請求項 1 記載の文書要約装置。

【請求項 3】 前記構成認定手段は、前記各位置の近傍領域における前記結束度を、各位置の前後に設定した 2 つの窓の中に含まれる語彙の類似性から求め、窓幅を段階的に縮小しながら前記話題境界の認定を繰り返すことを特徴とする請求項 2 記載の文書要約装置。

【請求項 4】 前記構成認定手段は、前記結束度を移動平均した値を、移動平均の開始点における右結束力および移動平均の終了点における左結束力として扱い、右結束力と左結束力が拮抗している位置の近傍を話題境界の候補区間と認定する候補区間認定手段を含み、該候補区間を用いて話題境界を認定することを特徴とする請求項 2 記載の文書要約装置。

【請求項 5】 前記結束度の比較的小さな領域を処理対象から除外して、該結束度の比較的大きな領域を重要部分として抽出する重要箇所特定手段をさらに備え、前記選択手段は、該重要部分に対応する話題のまとまりから前記重要文を選択することを特徴とする請求項 2 記載の文書要約装置。

【請求項 6】 前記抽出手段は、前記話題のまとまりの範囲に出現する語彙が、該話題のまとまりに対して特徴的であるかどうかを評価し、評価結果に基づいて前記重要語を抽出することを特徴とする請求項 1 記載の文書要約装置。

【請求項 7】 前記抽出手段は、評価対象の語彙が前記話題のまとまりに出現する頻度と、該評価対象の語彙が該話題のまとまりを含む大きな話題のまとまりに出現する頻度とを用いて、前記評価結果を得ることを特徴とする請求項 6 記載の文書要約装置。

【請求項 8】 前記抽出手段は、要約作成対象の話題のまとまりから局所的な重要語を抽出し、該要約作成対象の話題のまとまりを含む大きな話題のまとまりから大局的な重要語を抽出し、前記選択手段は、該局所的な重要語と大局的な重要語の両方の出現状況に基づいて、該要

約作成対象の話題のまとまりから前記重要文を選択することを特徴とする請求項 1 記載の文書要約装置。

【請求項 9】 前記要約の大きさに応じて、前記話題のまとまりに単語が出現する頻度と、該処理対象の話題のまとまりに単語が出現する頻度とをさらに備えることを特徴とする請求項 1 記載の文書要約装置。

【請求項 10】 与えられた文書中の処理対象の話題のまとまりに単語が出現する頻度と、該処理対象の話題のまとまりを含む大きな話題のまとまりに該単語が出現する頻度とを用いて、該単語が該処理対象の話題のまとまりに特徴的であるかどうかを評価し、評価結果に基づいて該処理対象の話題のまとまりから重要語を抽出する抽出手段と、

前記重要語の出現状況に応じて要約を生成する生成手段と、

前記要約を出力する出力手段とを備えることを特徴とする文書要約装置。

【請求項 11】 要約作成対象の話題のまとまりから局所的な重要語を抽出し、該要約作成対象の話題のまとまりを含む大きな話題のまとまりから大局的な重要語を抽出する抽出手段と、

前記局所的な重要語と大局的な重要語の両方の出現状況に基づいて、要約を生成する生成手段と、

前記要約を出力する出力手段とを備えることを特徴とする文書要約装置。

【請求項 12】 与えられた文書中の各位置の近傍領域における語彙的結束度を求める手段と、前記結束度の比較的小さな領域を処理対象から除外して、該結束度の比較的大きな領域を重要部分として抽出する重要箇所特定手段と、

前記重要部分を用いて要約を生成する生成手段と、前記要約を出力する出力手段とを備えることを特徴とする文書要約装置。

【請求項 13】 コンピュータのためのプログラムを記録した記録媒体であって、与えられた文書中の話題の階層的構成を認定するステップと、

認定された各話題に関する重要語を抽出するステップと、

前記重要語の出現状況に応じて、各話題のまとまりから重要文を選択するステップと、

前記重要文を用いて要約を生成するステップとを含む処理を前記コンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項 14】 コンピュータのためのプログラムを記録した記録媒体であって、

与えられた文書中の処理対象の話題のまとまりに単語が出現する頻度を求めるステップと、

前記処理対象の話題のまとまりを含む大きな話題のまとまりに前記単語が出現する頻度を求めるステップと、

得られた 2 つの頻度を用いて、前記単語が前記処理対象

の話題のまとまりに特徴的であるかどうかを評価するステップと、

評価結果に基づいて、前記処理対象の話題のまとまりから重要語を抽出するステップと、

前記重要語の出現状況に応じて要約を生成するステップとを含む処理を前記コンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項15】 コンピュータのためのプログラムを記録した記録媒体であって、

要約作成対象の話題のまとまりから局所的な重要語を抽出するステップと、

前記要約作成対象の話題のまとまりを含む大きな話題のまとまりから大局的な重要語を抽出するステップと、

前記局所的な重要語と大局的な重要語の両方の出現状況に基づいて、要約を生成するステップとを含む処理を前記コンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項16】 コンピュータのためのプログラムを記録した記録媒体であって、

与えられた文書中の各位置の近傍領域における語彙的結束度を求めるステップと、

前記結束度の比較的小さな領域を処理対象から除外して、該結束度の比較的大きな領域を重要部分として抽出するステップと、

前記重要部分を用いて要約を生成するステップとを含む処理を前記コンピュータに実行させるためのプログラムを記録したコンピュータ読み取り可能な記録媒体。

【請求項17】 与えられた文書中の話題の階層的構成を認定し、

認定された各話題に関する重要語を抽出し、

前記重要語の出現状況に応じて、各話題のまとまりから重要文を選択し、

前記重要文を用いて要約を生成することを特徴とする文書要約方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、自然言語などで書かれた機械可読文書の要約を行う装置およびその方法に関し、主として、長めのマニュアルや報告書などの要約(ダイジェスト)を作成し、文書の選別・閲覧のプロセスを支援することを意図している。

【0002】

【従来の技術】文書を要約するための主要な技術として、文書中の重要語を手掛かりに文を抽出(抜粋)して要約を作成する技術と、文書中の話題のまとまりを認定する技術がある。そこで、これらの従来技術について説明する。

【0003】まず、要約の作成技術について説明する。従来の文書の要約作成の技術には、大きく分けて2つの

方法がある。第1の方法は、文書において重要な部分を認定し、それを抜粋することで要約を作成する方法である。文書の重要な部分は、通常、節、段落、文などの論理要素の単位で抜粋される。以下では、これらを「文」という用語で代表させることにする。

【0004】第2の方法は、要約として抽出すべき情報の型紙を用意して、その型紙の条件にあった文書中の語句を抽出して要約としたり、その型紙によくあてはまる文を抽出して要約とする方法である。

【0005】第1の方法は、さらに、何を手掛かりに文の重要性を評価するかによっていくつかの方法に分類される。代表的な方法としては、次の3つが挙げられる。

(1) 文書中に出現する単語の頻度と分布を手掛かりとする方法

(2) 文と文とのつながりや文の出現位置を手掛かりとする方法

(3) 文の構文的パターンによって重要性を評価する方法

これらのうち、(1)の方法は、まず、文書に含まれる単語(語句)の重要度を決定し、次に、重要な単語をどれ位含んでいるかによって文の重要度を評価する。そして、評価結果に基づいて重要な文を選択して要約を作成する。

【0006】単語の重要度を決定する方法としては、文書中の単語の出現頻度(出現度数)そのものを用いる方法、単語の出現度数と一般的な文書集合におけるその単語の出現度数とのずれなどを加味して重みを付ける方法、単語の出現位置に応じて重みを付ける方法などが知られている。単語の出現位置に応じて重みを付ける場合は、例えば、見出しに出現する語を重要とみなすなどの処理が付加される。

【0007】ここで、対象とする単語は、日本語であれば自立語(特に名詞)のみに、英語であれば内容語のみに限るのが通例である。自立語・内容語とは、実質的な意味を持つ名詞、形容詞、動詞などの語であり、助詞、前置詞、形式名詞など、専ら構文的役割を果たすために使われる語とは区別される。なお、日本語の自立語の形式的定義は、独立した文節を構成できる語というものであるが、ここでは、上述の区別により自立語を定義している。

【0008】このような要約作成方法には、例えば、次のようなものがある。特開平6-259424「文書表示装置及び文書要約装置並びにデジタル複写装置」とその発明者による文献(亀田雅之、擬似キーワード相関法による重要キーワードと重要文の抽出、言語処理学会第2回年次大会発表論文集、pp. 97-100、1996年3月。)では、見出しに含まれる単語を多く含む部分を、見出しに関連の深い重要部分として抜粋することで要約を作成している。

【0009】特開平7-36896「文書を要約する方

法および装置」では、文書中に現れる表現（単語など）の複雑さ（語の長さなど）から重要な表現の候補（シード）を選び、重要性の高いシードをより多く含む文を抜粋することで要約を作成している。

【0010】特開平8-297677「主題の要約を生成する自動的な方法」では、文書内の単語の出現頻度が大きい順に「主題の用語」を認定し、重要な「主題の用語」を多く含む文を抽出することで要約を作成している。

【0011】特開平2-254566「自動抄録生成装置」では、出現頻度の大きい単語を重要語と認定し、重要語が初めて登場する部分や、重要語を多く含む部分、自動的に認定した意味段落の先頭に出現している文などを抜粋することで要約を作成している。

【0012】次に、文書中の話題のまとまりを認定する方法について説明する。この方法には、大きく分けて次の2つが挙げられる。

(1) 文書中で繰り返される語による話題の意味的な結び付き（語彙的結束性：lexical cohesion）に基づく方法

$$\sin(b_1, b_r) = \frac{\sum_i w_{t, b_1} w_{t, b_r}}{\sqrt{\sum_i w_{t, b_1}^2 \sum_i w_{t, b_r}^2}} \quad (1)$$

【0015】ここで、 $b_1$  と  $b_r$  は、それぞれ、左窓（文書の冒頭側の窓）、右窓（文書の末尾側の窓）に含まれる文書の部分を表し、 $w_{t, b_1}$ 、 $w_{t, b_r}$  は、それぞれ、左窓、右窓に出現する単語  $t$  の出現頻度を表す。また、(1) 式の右辺の  $\sum_i$  は、単語  $t$  に関する総和を表す。

【0016】(1) 式の類似度は、左右の窓に含まれる語彙に共通のものが多きほど大きくなり（最大1）、共通のものが少ない時に0となる。つまり、この値が大きい部分は、左右の窓で共通の話題を扱っている可能性が高く、逆に、この値が小さい部分は、話題の境界である可能性が高いことになる。

【0017】Hearst法は、(1) 式の値を文書の冒頭から  $d_s = (C_p - C_{mp}) + (C_r - C_{mp})$

そして、 $d_s$  が次式のような閾値  $h$  を越えた極小点だけを話題境界として認定している。

$$h = C_0 - \sigma / 2$$

ここで、 $C_0$ 、 $\sigma$  は、それぞれ、文書全体における類似度の平均値と標準偏差である。この方法によれば、類似度がより大きく落ち込んだ部分ほど、話題境界である可能性がより高いとみなされる。また、Hearstは、別法として、繰り返し出現する語の連鎖の開始・終了を手掛かりとして、開始点・終了点の近傍に話題境界を認定する方法も示している。

【0021】語彙的結束性に基づいて話題のまとまりを認定する方法としては、その他に、日本語の提題助詞

(2) 接続詞などで示される文間の接続関係 (coherence relation) から文章構造 (rhetorical structure) を求める方法

これらのうち、(1) の語彙的結束性に基づく方法として、まず、Hearstの方法 (Marti A. Hearst, Multi-paragraph segmentation of expository text, In Proceedings of the 32nd Annual Meeting of Association for Computational Linguistics, pp. 9-16, 1994.) を簡単に説明する。

【0013】この方法（以下、Hearst法と称する）は、意味的に関係の深い部分には、同一の語彙が繰り返し出現するという性質（語彙的結束性）を利用して、話題の切れ目となる部分を自動的に認定するものである。この方法では、まず、文書中の各位置の前後に、段落程度の大きさ（120語程度）の窓を設定し、その2つの窓にどれくらい同じ語彙が出現しているかを表す類似度を測定する。類似度としては、次式のような余弦測度 (cosine measure) と呼ばれる値が用いられている。

【0014】

【教1】

ら末尾まである間隔（20語）で測定し、極小となる位置を話題境界と認定するものである。このとき、類似度の細かい振動を無視するために、次のような調整を行っている。まず、極小点  $m$  の周囲の部分を取り出す。この部分には、極小点の左側の単調減少している部分と極小点の右側の単調増加している部分が含まれる。

【0018】次に、切り出された部分の開始点  $l_p$ 、終了点  $r_p$  における類似度  $C_{lp}$ 、 $C_{rp}$  と、類似度の極小値  $C_m$  との差をもとに、次式の値  $d_s$  (depth score) を計算し、これを極小点における類似度の変動量の指標とする。

【0019】

$$d_s = (C_{lp} - C_m) + (C_{rp} - C_m) \quad (2)$$

【0020】

$$(3)$$

「は」のついた文節で始まる文（例えば、「Hearstは、」で始まる文）などを手掛かりとする方法が知られている（特開平7-160711「書き言葉テキストに対する話題構造認識方法および装置」）。また、その方法とHearstの別法に類似する方法とを併用する方法も知られている（望月源、本田岳夫、奥村孝、重回帰分析とクラスタ分析を用いたテキストセグメンテーション、言語処理学会第2回年次大会発表論文集、pp. 325-328、1996年3月。）。

## 【0022】

【発明が解決しようとする課題】しかしながら、上述した従来の文書要約方法には次のような問題がある。文書における重要語を決定し、重要語を多く含む文を抜粋することで文書の要約を作成する方法では、長めの文書、特に複数の話題に関する文章が混在している複合文書の要約を作成するのが困難である。複数の話題に関する文章が混在している場合、話題毎に重要な単語が異なる可能性が高いので、文書中出現頻度の大きい単語を単純に重要語とみなすことができない。単純に重要語を決定してしまつと、ある話題に関する重要性を手掛かりに、別の話題の部分から重要でない文が抜粋されてしまうことがあるからである。

【0023】この問題を解決するためには、文書中の話題のまとまりを認定する必要がある。ここで、語彙的結束性から直接に大きな話題のまとまりを認定する方法がないというもう1つの問題がある。

【0024】従来の技術では、語彙的結束性に基づいて話題のまとまりを認定する場合、Bears法のように、数段落程度の大きさあるいは大きくも新聞の1つの記事程度までの大きさのまとまりの認定しか試みられていなかった。そして、それより大きなまとまりは、原文書の章などの書式を手掛かりに認定していた。

【0025】例えば、前述の特開平2-254566では、内容的に関連度の高い一連の形式段落（字下げなどにより形式的に区切られた段落）を意味段落として自動認定し、文書全体で出現頻度の大きい語だけでなく、それぞれの意味段落で出現頻度の大きい語も重要語として抽出して、要約を作成している。しかし、この方法では、書式を手掛かりに認定した章や節などの切れ目を、自動的に認定した意味段落の分割点より優先しているため、意味段落は章や節などの切れ目を越えることがなく、より大きな話題のまとまりの認定は試みられていない。

【0026】話題を認定する方法としても、意味段落認定の主要手掛かりは、形式段落2つ分の範囲で繰り返される語彙であるので、大きな話題のまとまりを認定することは困難である。また、語彙が初めて出現する位置の情報も用いているが、大きな間隔で繰り返される語による結束性などを判定するには十分とは言えない。

【0027】しかしながら、同じ章に属する節であっても意味的なまとまり方に違いのある場合もあり、このような場合にも的確に大きな話題のまとまりを認定する方法が必要である。また、文書の書式などは特定の種類の文書に関する約束事であるため、様々な種類の文書の要約を行うためには、文書の種類毎にいちいち経験的な規則を用意しなければならないという問題もある。

【0028】本発明の課題は、語彙的結束性のような文章一般に見られる現象をもとに、文書中の話題構成を自動的に認定して、話題構成に対応する要約を作成する汎

用の文書要約装置およびその方法を提供することである。

## 【0029】

【課題を解決するための手段】図1は、本発明の文書要約装置の原理図である。図1の文書要約装置は、構成認定手段1、抽出手段2、選択手段3、および出力手段4を備える。

【0030】構成認定手段1は、与えられた文書中の話題の階層的構成を認定し、抽出手段2は、認定された各話題に関する重要語を抽出する。また、選択手段3は、重要語の出現状況に応じて、各話題のまとまりから重要文を選択し、重要文を用いて要約を生成する。出力手段4は、生成された要約を出力する。

【0031】ここで、話題の階層的構成とは、文書を構成する複数の話題のまとまりが2段以上の階層構造を成していることを意味する。この階層的構成は、例えば、文書を構成する複数の大きな話題のまとまりの各々が、1つ以上のより小さな話題のまとまりを含み、小さな話題のまとまりの各々が、1つ以上のさらに小さな話題のまとまりを含むというような話題の包含関係に対応する。

【0032】構成認定手段1は、例えば、文書全体の大きさの $1/4 \sim 1/10$ 程度から段落程度の大きさまで、数種類の大さきの窓幅を設定し、語彙的結束性の強さを表す結束度を各窓幅で測定する。これにより、大きな間隔で繰り返される語などによる大局的な結束性と、小さな間隔で繰り返される語などによる局所的な結束性の両方を捉えることができ、大きな話題のまとまりから小さな話題のまとまりに至る話題の階層的構成を認定することができる。

【0033】抽出手段2は、例えば、処理対象の話題のまとまりに単語が出現する頻度と、その話題のまとまりを含む大きな話題のまとまりにその単語が出現する頻度とを用いて、その単語が処理対象の話題のまとまりに特徴的であるかどうかを評価する。そして、その評価結果に基づいて、処理対象の話題のまとまりから重要語を抽出する。

【0034】このように、処理対象の話題のまとまりから重要語を抽出する際に、それを含む上位の話題のまとまりも参照するため、その単語の重要性を上位の話題のまとまりとの関係から評価することができる。このため、話題に関わらず単に多く出現する語を誤って重要語と判定することなく、効率的に重要語を抽出できる。

【0035】また、抽出手段2は、例えば、要約作成対象の話題のまとまりから局所的な重要語を抽出し、その話題のまとまりを含む大きな話題のまとまりから大局的な重要語を抽出する。そして、選択手段3は、局所的な重要語と大局的な重要語の両方の出現状況に基づいて、要約作成対象の話題のまとまりから重要文を選択し、要約を生成する。

【0036】このように、要約作成対象の話題のまとまりから重要文を選択する際に、それを含む上位の話題のまとまりに出現する重要語も参照するため、局所的な話題に関する文と大局的な話題に関する文の両方をバランスよく含んだ要約を生成することができる。

【0037】例えば、図1の構成認定手段1は、後述する図2の話題構成認定部26に対応し、図1の抽出手段2は図2の重要語抽出部29に対応し、図1の選択手段3は図2の重要文選択部30に対応し、図1の出力手段4は図2の出力部31に対応する。

【0038】

【発明の実施の形態】以下、図面を参照しながら、本発明の実施の形態を詳細に説明する。図2は、本発明の文書要約装置の基本構成を示している。図2において、文書要約装置 (digest generator) 12は、要約対象文書 (input document) 11が入力されると、その要約文書13を作成して出力する。

【0039】文書要約装置12は、入力部 (input unit) 21、単語認定部 (tokenizer) 22、単語辞書 (machine readable dictionary) 24、要約粒度決定部25、話題構成認定部 (topic structure detector) 26、重要箇所特定部28、重要語抽出部 (keyword extractor) 29、重要文選択部 (sentence selector) 30、および出力部31を備える。

【0040】入力部21は、要約対象文書11を読み込み、単語認定部22に渡す。単語認定部22は、形態素解析部 (morphological analyzer) 23を含み、それを用いて要約対象文書11を言語的に解析して、文書11に含まれる内容語 (名詞・動詞・形容詞・形容動詞など) を切り出す。このとき、形態素解析部23は、単語辞書24を参照して、文書11中の文を、品詞情報付きの単語リストに変換する。単語辞書24は、形態素解析用の単語辞書であって、単語の表記文字列と品詞・活用の情報との対応関係などを記述している。

【0041】要約粒度決定部25は、文書11の大きさと望ましい要約の大きさから、要約として抽出すべき話題の数を計算し、要約を作成する単位とする話題のまとまりの大きさを決定する。

【0042】話題構成認定部26は、話題境界候補区間認定部 (topic boundary detector) 27を含み、それを用いて共通の話題について記述している文書の部分 (話題のまとまり) を自動認定する。話題境界候補区間認定部27は、話題構成認定部26のサブモジュールとして、語彙的結束度の小さい区間を話題境界の候補区間として認定する。語彙的結束度とは、文書11中の各位置の近傍領域における語彙的結束性の強さを表す指標であり、例えば、各位置の前後に設定したある幅の窓内に出現する語彙の類似性から求められる。

【0043】重要箇所特定部28は、語彙的結束度の小さい話題のまとまりを以後の処理の対象から除外し、文

書の主要な部分だけが要約に出力されるようにする。重要語抽出部29は、話題構成認定部26が認定した話題のまとまりについて、その範囲に出現する語彙が話題に特徴的であるかどうかを評価し、特徴的に出現している語を重要語として抽出する。

【0044】重要文選択部30は、それぞれの話題のまとまりについて、重要語を多く含む文を選択し、選択した文を原文書11での出現順に並べる。そして、必要に応じて選択されなかった文の存在を表す印や段落境界などを挿入することで、要約文書13を作成する。出力部31は、作成された要約文書13を処理結果として出力する。

【0045】図2の文書要約装置12によれば、話題構成認定部26が、共通の話題について記述している文書の部分を話題のまとまりとして認定し、重要語抽出部29が、それぞれの話題のまとまりに特徴的な語を抽出する。このため、話題の異なる複数の文章が混在している複合文書に対しても、精度よく重要語を抽出することができる。また、重要文選択部30が、話題のまとまり毎に、そのまとまりに特徴的な重要語を手掛かりに重要文を選択して要約を作成するので、別の話題の重要語の影響で不要な文が抜かれてしまうこともない。

【0046】話題構成認定部26は、語彙的結束度に基づき話題を認定する際に、文書全体の大きさの $1/4 \sim 1/10$ 程度の大きさの窓により測定した語彙的結束度から、段落程度の大きさの小さい窓により測定した語彙的結束度まで、数種類の語彙的結束度を併用する。このように、大きな間隔で繰り返される語などによる大局的な結束性と、小さな間隔で繰り返される語などによる局所的な結束性の両方に関する情報を利用しているので、話題構成認定部26は、大きな話題のまとまりから小さな話題のまとまりまで、もれなく話題のまとまりを認定できる。

【0047】さらに、話題構成認定部26のサブモジュールである話題境界候補区間認定部27は、各窓幅の語彙的結束度を移動平均した値を、移動平均の開始点における右結束力および移動平均の終了点における左結束力として扱い、右結束力と左結束力が拮抗している部分 (結束力拮抗点) の近傍を、話題境界の候補区間と認定する。

【0048】移動平均を用いることで、語彙的結束度の小さな変動、すなわち、移動平均区間 (移動平均をとる区間) の大きさに比べて小さな範囲内の変動が平滑化される。このため、それぞれの結束力拮抗点の間隔は、ほとんどが移動平均区間の大きさ程度以上となる。これにより、話題構成認定部26は、それぞれの窓幅の語彙的結束度に基づいて、窓幅程度 (移動平均区間の幅程度以上) の大きさの話題のまとまりを選択的に認定できるので、話題の階層的構成を正確に認定することができる。

【0049】また、重要語抽出部29は、統計的検定法

により、それぞれの話題のまとまりにおいて有意に多く現れると判定された語を重要語と認定する。このため、話題に関わらず単に多く出現する語を誤って重要語と判定することなく、効率的に重要語を抽出できる。

【0050】さらに、重要語抽出部29を使って、要約作成対象の話題のまとまりからだけでなく、要約作成対象の話題のまとまりを含むより大きな話題のまとまりからも大域的な重要語を抽出することができる。これにより、小さな話題のまとまりが並んで、より大きな話題のまとまりを構成しているような場合にも、適切に重要語を抽出することができる。すなわち、個々の小さな話題に特徴的な重要語（副主題を表す語）と、それらに共通する大きな話題に特徴的な重要語（主題を表す語）の両方を区別して抽出することができる。

【0051】そして、重要文選択部30は、主題を表す語と副主題を表す語の両方を手掛かりに重要文を選択して要約を作成する。このため、主題と副主題の両方をバランスよく含んだ要約が作成される。

【0052】また、話題構成認定部26は、要約粒度決定部25が決定した大きさ程度の話題のまとまりを認定し、重要語抽出部29と重要文選択部30が、この話題のまとまりを単位に要約を作成するので、結果として、抽出すべき話題の数程度で、かつ、同じ程度の大きさの話題を、バランスよく要約に取り込むことができる。

【0053】さらに、重要箇所特定部28は、話題構成認定部26が認定した話題のまとまりのうち、結束度の小さい区間を要約対象から除外する。このため、単に項目を列挙しただけの部分などを按捺してしまうことがなく、内容の濃い要約を作成することができる。

【0054】図2の文書要約装置12は、例えば、図3に示すような情報処理装置（コンピュータ）を用いて構成することができる。図3の情報処理装置は、出力装置41、入力装置42、CPU（中央処理装置）43、ネットワーク接続装置44、媒体駆動装置45、補助記憶装置46、およびメモリ（主記憶）47を備え、それらはバス48により互いに接続されている。

【0055】メモリ47は、例えば、ROM（read only memory）、RAM（random access memory）などを含み、文書要約処理に用いられるプログラムとデータを格納する。ここでは、図2に示した入力部21、単語認定部22、形態素解析部23、要約粒度決定部25、話題構成認定部26、話題境界候補区間認定部27、重要箇所特定部28、重要語抽出部29、重要文選択部30、および出力部31が、プログラムモジュールとして格納されている。CPU43は、メモリ47を利用してプログラムを実行することにより、必要な処理を行う。

【0056】出力装置41は、例えば、ディスプレイやプリンタなどであり、ユーザへの問い合わせや要約文書13などの出力に用いられる。入力装置42は、例えば、キーボード、ポインティングデバイス、タッチパネ

ルなどであり、ユーザからの指示や要約対象文書11の入力に用いられる。

【0057】補助記憶装置46は、例えば、磁気ディスク装置、光ディスク装置、光磁気ディスク（magneto-optical disk）装置などであり、要約対象文書11、要約文書13、単語辞書24などの情報を格納する。この補助記憶装置46に、上述のプログラムとデータを保存しておき、必要に応じて、それらをメモリ47にロードして使用することもできる。

【0058】媒体駆動装置45は、可搬記録媒体49を駆動し、その記録内容にアクセスする。可搬記録媒体49としては、メモリカード、フロッピーディスク、CD-ROM（compact disk read only memory）、光ディスク、光磁気ディスクなど、任意のコンピュータ読み取り可能な記録媒体が用いられる。この可搬記録媒体49に上述のプログラムとデータを格納しておき、必要に応じて、それらをメモリ47にロードして使用することもできる。

【0059】ネットワーク接続装置44は、LAN（local area network）などの任意のネットワーク（回線）を介して外部の装置と通信し、通信に伴うデータ変換を行う。また、必要に応じて、上述のプログラムとデータを外部の装置から受け取り、それらをメモリ47にロードして使用することもできる。

【0060】図4は、図3の情報処理装置にプログラムとデータを供給することのできるコンピュータ読み取り可能な記録媒体を示している。可搬記録媒体49や外部のデータベース50に保存されたプログラムとデータは、メモリ47にロードされる。そして、CPU43は、そのデータを用いてそのプログラムを実行し、必要な処理を行う。

【0061】次に、図2の文書要約装置12の各モジュールの動作を、具体例を用いてより詳細に説明する。要約対象文書としては、（社）電子工業振興協会「自然言語処理システムの動向に関する調査報告書」（平成9年3月）第4章「ネットワークアクセス技術専門委員会活動報告」（pp. 117-197）を用いた。以下の実施形態では、この文書から文を按捺してA4、1〜2枚（1500文字）程度の要約の作成を試みる。

【0062】従来、要約の大きさとしては、原文書の1/4程度の大きさが目安とされてきたが、この要約対象文書は81ページの大きさを持ち、従来の自動要約技術が対象としてきた新聞の社説や記事、数頁程度の論文などに比べて巨大である。また、オンラインで文書を閲覧する場合、画面に一度に表示できるのは2ページ程度が限度である。これらの条件を考慮して、上述のような要約の大きさが決められている。

【0063】要約対象文書の全体を掲載することは適当ではないので、参考として、要約対象文書中の見出しの一覧を図5から図7に示す。図5は、4.1節および

4. 2節の見出しを出現順に示しており、図6は、4. 3節の見出しを出現順に示しており、図7は、4. 4節の見出しを出現順に示している。

【0064】図8は、単語認定部22による単語認定処理のフローチャートである。単語認定部22は、まず、要約対象文書に形態素解析を施し、品詞付きの単語リストを作成する(ステップS11)。次に、品詞を手掛かりに内容語(名詞・動詞・形容詞・形容動詞)を認定し、内容語に対応する文書の部分に印を付けて(ステップS12)、処理を終了する。図9は、要約対象文書の冒頭部分を示しており、図10は、単語認定部22からの対応する出力を示している。

【0065】図8のステップS11において、形態素解析部23は、図11に示すような形態素解析処理を行う。形態素解析部23は、まず、単語リストをクリアし(ステップS21)、文書の先頭から句点(またはピリオド)などを手掛かりに文の取り出しを試み(ステップS22)、文が取り出せたかどうかを判定する(ステップS23)。

【0066】文が取り出せれば、次に、単語辞書24を参照して、文に含まれている単語の候補を求める(ステップS24)。日本語の場合は、図9に示したように、単語と単語の境界が形式的に明示されていないので、文に含まれる部分文字列に対応するすべての単語を候補として求める。例えば、「東京都は大都市だ」という文が取り出された場合、図12に示すように、この文に含まれるすべての部分文字列が単語の候補となる。

【0067】これに対して、英語の場合は、単語の境界が空白(スペース)により明示されているため、空白で区切られた文字列に対応する単語について、品詞の候補を求めることが主な処理となる。例えば、「Tokyo is the Japanese capital.」という文が取り出された場合、図13に示すように、この文に明示的に含まれる5つの単語の基本形と品詞が求められる。

【0068】次に、形態素解析部23は、品詞レベルの接続の観点から、妥当な単語の並びを選択し(ステップS25)、選択された単語の並びに品詞と出現位置の情報を付加して、出現順に単語リストに追加する(ステップS26)。次に、次の文の取り出しを試み(ステップS27)、ステップS23以降の処理を繰り返す。そして、ステップS23において文が取り出せなくなると、処理を終了する。

【0069】図10の単語認定結果において、墨付き括弧で括弧された部分が形態素解析部23の認定した内容語である。内容語が活用語(動詞・形容詞)の場合、墨付き括弧内で、スラッシュ(/)の前の部分は語幹を表し、スラッシュの後の部分は終止形の活用語尾を表す。これは、後の処理で単語の区別を行うために用いられる情報であるが、この情報の代わりに、品詞と活用を付加

しておいてもよい。要するに、例えば、「い/る」と「い/く」のように、語幹だけでは区別の付かない単語を区別するための識別情報であれば、任意のものを用いることができる。

【0070】また、ステップS25において、単語の並びの妥当性を評価する方法は、形態素解析法として各種のものが知られており、任意のものを用いることができる。例えば、単語の並びの妥当性を訓練データにより推定された出現確率を用いて評価する方法が報告されている(Eugene Charniak, Hidden markov models and two applications. In Statistical Language Learning, chapter 3, pp.37-73. The MIT Press, 1993. / Masaaki Nagata, A stochastic Japanese morphological analyzer using a forward-DP backward-N-best search algorithm. In Proceedings of COLING '94, pp.201-207, 1994. / 永田昌明, 前向きDP後向きA\*アルゴリズムを用いた確率的日本語形態素解析システム, 情報処理学会, 1994年5月。)

【0071】なお、図10の例では、単語認定部22がすべての内容語を切り出しているが、切り出しの対象を名詞だけに絞っても構わない。また、英語の文書を対象に処理する場合には、形態素解析処理を行う代わり、空白で区切られたすべての語のうち、話題に関わらずどこにでも出現する語彙(冠詞、前置詞などの機能語や特に高い頻度で出現する語)を取り除いて、単語を切り出してよい。このような処理は、単語辞書24の代わりに、機能語や特に高い頻度で出現する語を格納したストップワードリスト(stop word list)を用意すれば、容易に実現できる。

【0072】図14は、要約粒度決定部25による要約粒度決定処理のフローチャートである。要約粒度決定部25は、まず、望ましい要約の大きさ $S_a$ 、望ましい各話題の抜粋量 $S_c$ 、最小窓幅 $w_{\min}$ 、窓幅 $b$ の4つのパラメータをユーザから受け取り(ステップS31)、 $S_a$ を $S_c$ で割って抽出すべき話題の個数 $N_c$ を求める(ステップS32)。

【0073】ところで、図14では、図面の見やすさを考慮して、記号“ $w_{\min}$ ”の添字を、“ $w_{\min}$ ”のように下線を付加して記している。他の添字についても、同様の表記法が用いられる場合がある。

【0074】次に、要約粒度決定部25は、要約対象文書の延べ語数 $W_d$ を求める(ステップS33)。そして、 $W_d$ を $N_c$ で割って抽出すべき話題の大きさの目安 $w_c$ を算出した後、初項を $w_{\min}$ とし公比を $r$ とする等比級数の中から、 $w_c$ を超えない最大の項を選んで基本窓幅 $w_i$ とし(ステップS34)、処理を終了する。このとき、 $w_i$ は、次式により計算される。

【0075】

$$w_i = w_{\min} * (r * \text{rint}(\log_r((W_d / N_c) / w_{\min})))$$



ここで、 $**r$ は $r$ を $\text{int}()$ 乗することを表し、 $\text{int}()$ は、括弧内の部分の少数点以下を切り捨てて、整数にすることを表している。等比級数の各項は、以降の処理で話題の階層的構成を認定する際に、各層における処理用の窓幅として用いられる。

【0076】また、別法として、 $\text{int}(W_d / N_i)$ の値をそのまま $w_i$ として使い、 $w_{\text{min}}$ の方を $w_i * (1/r) ** n$  ( $n$ は整数) の形で定義してもよい。さらに、等比級数によらずに、他の任意の方法で窓幅を段階的に縮小していくことも可能である。ただし、後述するように、 $w_{\text{min}}$ を固定し、2の累乗の値を公比とする等比級数を用いる方法が、計算効率上望ましいことが分かっている。

【0077】例えば、 $S_s = 1500$  (文字)、 $S_t = 150$  (文字)、 $w_{\text{min}} = 40$  (語)、 $r = 2$ 、 $W_d = 17816$  (語) とすると、抽出すべき話題の数 $N_t$ は $10(1500\text{文字}/150\text{文字})$ となる。この場合、話題のまとまりの大きさの目安 $w_t$ は $1800$ 語程度 ( $17816\text{語}/10$ ) になるので、これを越えない値として、 $1280$ 語 ( $40 * 2^5$ ) が基本窓幅 $w_i$ に採用される。

【0078】新聞記事などの要約実験によれば、話題の内容を理解可能にするためには、それぞれの話題に関して3文程度 (見出し1文+2~3文:  $120 \sim 150$ 文字程度) 以上の文を採捨する必要があるという経験的知識が得られている。上記の採捨量 $S_t$ の値は、このような経験的知識に基づき決定されている。また、窓幅 $w_{\text{min}}$ の値は、新聞記事やレポートなどの平均的な語数により決定されている。

【0079】次に、話題構成認定部26の処理について説明する。本実施形態においては、話題のまとまりは前述のHearst法を拡張した方法により認定される。したがって、文書の各位置の語義的結束度 (以下、結束度と略す) を測定し、結束度の小さい部分に話題境界を認定するという方針が採用されている。本実施形態とHearst法の主な違いは、次のような点にある。

(1) 結束度を測定するための窓幅の違い

本実施形態では、結束度の計算に用いる窓として、Hearst法より巨大なもの (要約対象文書の全体の語数の $1/4 \sim 1/10$ 程度: 上述の例では $1280$ 語) から、段落程度の大きさ (数十語から $100$ 語程度: 上述の例では $40$ 語) のものまで、幅の異なるものを数種類併用している。

(2) 話題境界の認定手順および認定対象とする話題境界の違い

本実施形態では、Hearst法のように、異なる窓幅で測定したそれぞれの結束度 (または類似度) について、結束度が極小となる位置をそのまま話題境界と認定するのではなく、結束度の移動平均 (moving average) を用いる

## (4)

ことで、窓幅程度の大きさのまとまりのみを話題のまとまりとして認定している。

【0080】これらの違いは、本発明が文書中の話題の階層的構成を認定することに起因する。ここで、話題の階層的構成とは、例えば、ある話題を扱った章の中に含まれるいくつかの小さな話題の節が含まれるような、話題の包含関係を有する構成のことである。

【0081】話題の階層的構成を認定する理由は、小さな話題の部分が並んで、より大きな話題のまとまりを形成している場合に、個々の小さな話題に特徴的な重要語 (副主題を表す語) と、それらに共通する大きな話題に特徴的な重要語 (主題を表す語) の両方を区別して抽出することで、主題と副主題の両方をバランスよく含んだ要約を作成するためである。

【0082】従来の研究では、数千語レベルの窓幅を使って測定した類似度が、実際の文章における話題の推移と対応するかどうか、すなわち、実際の文章の話題のまとまりの認定に使えるのかどうかは確かめられていなかった。これを確かめようとした研究がなかったのは、このような単純な測定法で数千語レベルの窓幅を使ってしまつと、測定結果が雑音だらけになって、無意味な変動しか示さないだろうという先入観があったものと推察される。

【0083】例えば、前述のHearstの文献の結論によれば、シンテラス (類語辞典) などのもつと複雑な情報を使ってより精密な境界認定を実現する可能性が示唆されているが、窓幅については、実験方法の説明の中で簡単に述べられているだけである。したがって、本実施形態のように、窓幅を極端に変更した場合にどうなるのかなどについての考察は見られない。

【0084】Hearstは、実験対象毎に微調整して提示した程度の窓幅が、この方法における最適値であり、文章中の副主題に関する数段落程度の大きさの話題のまとまり (passage) を認定するという問題設定が、この方法の限界であると考えていた可能性が高い。また、Hearstの目的は、このような数段落程度のまとまりを認定することに限定していたとも考えられる。

【0085】そこで、Hearstの文献で用いられているものより $5 \sim 10$ 倍程度巨大な窓幅によって測定した類似度が、意味のある変動を示すのかどうかを確かめるために、上述の要約対象文書の話題境界をHearst法により認定する実験を行った。この実験により得られた類似度を結束度としてプロットした結果、図15および図16のような結束度分布が得られた。

【0086】これらの図において、横軸の文書における位置は、文書の冒頭から各位置までの間に出現した内容語の延べ数を表す。また、点線は、要約対象文書内の各節の開始位置を表し、長い点線ほど大きい節に対応している。さらに、記号◇でプロットした折れ線グラフは、

(1) 式の余弦測度により求めた結束度の系列を表し、記号\*の付いた棒グラフは、結束度の極小点における(2)式のdepth scoreを表し、水平線は(3)式の閾値を表す。

【0087】結束度の計算においては、図15では1280語幅の窓が用いられ、図16では640語幅の窓が用いられている。また、結束度の系列は、それぞれの窓幅の1/8(160語または80語)の刻み幅で計算され、プロットされている。

【0088】図15および図16を見ると、点線で示された各節の開始位置付近に、閾値を越えるdepth scoreが付与されており、数千語レベルの窓幅を使って測定した結束度も意味のある変動を示すことが分かる。このように、巨大な窓幅による語彙的結束度を用いれば、章・節レベルの話題境界も認定可能である。また、図15および図16を比較すると、大きな窓幅の結束度を使うと大きな話題の切れ目が認定でき、小さな窓幅の結束度を使うと小さな話題の切れ目が認定できるという傾向も見てとれる。

【0089】しかしながら、この実験結果によれば、Hearst法における次のような問題点が指摘される。

(1) 大きな窓幅で認定した話題境界と小さな窓幅で認定した話題境界の対応付けが明確でない。

(2) 結束度がおおむね単調減少または単調増加している部分の途中に小さな極値が挟まるだけで、depth scoreは大きく変化してしまうので、これは安定な指標とは言えない。

【0090】これらの問題点は、例えば、要約対象文書の4.3節の末尾にある参考文献の部分(参)から4.4.1(1)節の部分までに対応する処理結果に現れている。図15では、この部分は、大局的に見れば、結束度の1つの谷である。この傾向は、図16でも変わらない。

【0091】しかし、図16では、4.3節(参)の部分の幅の狭い小さな山P1と、4.4.1(2)節の半ばから4.4.1(3)節まであたりの谷P2とが明確に現れている。このため、640語幅の話題境界は、1280語幅の話題境界と大きくずれており、このずれは図15の刻み幅以上に達する。

【0092】話題の階層的構成を認定する場合、4.4節の開始位置を大きな話題の切れ目と認定し、4.3節(参)の開始位置なにより小さな話題の切れ目と認定したい。しかし、Hearst法のdepth scoreは安定でないため、これを話題境界に対応する話題の大きさの指標とすることは無理がある。

【0093】また、depth scoreが安定でないため、大きな窓幅の結束度により認定された話題境界が、必ずしも小さな窓幅の結束度により認定されるとは限らない。さらに、大きな窓幅の結束度より大きな話題の切れ目だけが話題境界と認定されるわけでもない。このため、

Hearst法は、大きな窓幅で大きな話題の切れ目を認定し、小さな窓幅で小さな話題の切れ目を認定するという処理には使えない。

【0094】本実施形態の話題境界候補区間認定部27は、これらの問題点を解決するために、移動平均法を応用して、話題境界の区間推定を行う。このとき、各窓幅毎に、結束度の移動平均した値を、移動平均の開始点における右結束力および移動平均の終了点における左結束力として扱い、右結束力と左結束力の拮抗点の近傍を話題境界の候補区間と認定する。

【0095】移動平均を用いることで、結束度の小さな変動、すなわち、移動平均区間の大きさに比べて小さな範囲内の変動が平滑化される。このため、それぞれの結束力拮抗点の間隔は、ほとんどが移動平均区間の大きさ程度以上となる。

【0096】これにより、話題構成認定部26において、次のような話題の階層的構成の認定手順が実現される。

(1) 大きな窓幅では大きな話題に対応する話題境界だけを選択的に認定する。

(2) 話題境界を区間推定し、大きな窓幅による話題境界と、区間の範囲内で一致しているとみなせるより小さな窓幅による話題境界を求める。そして、両者を同一の話題境界とみなす。

【0097】図17は、話題構成認定部26による話題構成認定処理のフローチャートである。話題構成認定部26は、まず、基本窓幅 $w_1$ 、最小窓幅 $w_{min}$ 、窓幅比 $r$ の3つのパラメータを要約粒度決定部25から受け取り(ステップS41)、結束度を測定するための窓幅の集合 $W$ を求める(ステップS42)。窓幅の集合 $W$ は、初項 $w_0$ を $w_1 * r$ とし、公比を $1/r$ とする等比級数から、 $w_{min}$ 以上の大きさの項を集めて作成される。このとき、 $W$ における最大窓幅は、 $w_0 = w_1 * r$ となる。

【0098】なお、前述したように、窓幅の集合 $W$ の選び方はこれ以外の方法であってもよいが、計算効率上は、 $w_{min}$ を固定し、 $r$ として2の累乗の乗を用いる方法が望ましい。

【0099】 $w_1 = 1280$ (語)、 $w_{min} = 40$ (語)、窓幅比 $r = 2$ の場合は、最大窓幅 $w_0$ は2560語(1280\*2)となる。次に、話題構成認定部26は、図10に示したように、内容語に印が付けられた文書をもとに、文書中の各位置の結束度を、 $W$ 中のそれぞれの窓幅毎に計算し、結束度系列として記録する(ステップS43)。

【0100】ここでは、まず、文書の各位置(基準点)の前後に設定した2つの窓の中に出現している語彙(ここでは内容語)を比較し、共通している語彙が多い程大きくなるような値を計算して、その位置における結束度とする。そして、窓の位置を文書の冒頭から末尾に向か

って一定の刻み幅  $t_{ic}$  でずらしながら、結束度の計算を繰り返し、計算した結束度を、文書の冒頭から末尾に向かう系列として記録する。

【0101】なお、刻み幅  $t_{ic}$  は、窓幅より小さければいずれの値でも構わないが、例えば、窓幅の  $1/8$  というように、窓幅に比例するように設定するのが効率的である。この  $t_{ic}$  の値は、ユーザにより指定される。

【0102】図18は、図10の単語認定結果において設定された2つの窓を示している。ここでは、40番目の内容語「サービス/する」と41番目の内容語「左窓W1と右窓W2が設定されている。この位置における結束度は、次のように計算される。

右結束度＝共通語彙数/左窓中の出現語彙数

$$= 6/29$$

$$= 0.207$$

左結束度＝共通語彙数/右窓中の出現語彙数

$$= 6/29$$

$$= 0.207$$

結束度＝(右結束度+左結束度)/2

$$= 0.207$$

(5)、(6)、(7)式により得られた各結束度には、次のような意味がある。ある窓に含まれる語がその右側(文書の末尾へ向かう方向)の部分にも出現している場合、その数が多い程、その窓の部分は右側との結び付きが強いと考えられる。この指標が、(5)式の右結束度である。同様に、ある窓とその左側(文書の冒頭へ向かう方向)部分との結び付きの強さを示す指標が、(6)式の左結束度である。そして、基準点においてこれらの2種類の結び付きの強さを平均したものが、(7)式の結束度である。

(5)、(6)、(7)式により得られた各結束度には、次のような意味がある。ある窓に含まれる語がその右側(文書の末尾へ向かう方向)の部分にも出現している場合、その数が多い程、その窓の部分は右側との結び付きが強いと考えられる。この指標が、(5)式の右結束度である。同様に、ある窓とその左側(文書の冒頭へ向かう方向)部分との結び付きの強さを示す指標が、(6)式の左結束度である。そして、基準点においてこれらの2種類の結び付きの強さを平均したものが、(7)式の結束度である。

【0106】なお、結束度としては、(7)式の値でなくとも、文書中の各位置の近傍領域における語彙的結束性の強さを表す指標として妥当な値であれば、どんなものを用いてもよい。例えば、Bears法のように、左右の窓中の語彙の類似性を表す相対度測度を結束度として用いても構わない。

【0107】また、文書中の各位置の近傍領域を2つの窓に分割せずに、その近傍領域に一回以上出現している単語の数を結束度とすることもできる。実際、各位置を中心とする近傍領域に、類義語や関連語(例えば、「ウェ이터」と「レストラン」)などの意味的に関連する単語が出現する割合に対応する値を、結束度として用いることも報告されている(小嶋秀樹、古部延治、単語の結束性についてテキストを場面に分割する試み、電気情報通信学会、信学技報NLC93-7、1993年5月)。

【0108】ただし、(7)式に示した結束度の方が、計算が単純であり、解釈もしやすい。以下の説明において、(7)式の結束度を他の結束度と区別する必要があ

【0103】まず、図19に示すように、左窓W1と右窓W2中に出現している内容語の異なる数(窓中の出現語彙数)を数える。図18では、この出現語彙数は、W1、W2ともに29語ずつである。次に、左窓W1と右窓W2の両方に出現している内容語の異なる数(共通語彙数)を数える。図18では、W1、W2中に下線を付けて示した6語が共通語彙数となる。

【0104】最後に、左窓W1における共通語彙数と出現語彙数の比を右結束度とし、右窓W2における共通語彙数と出現語彙数の比を左結束度として、これらの結束度の算術平均を求め、これを基準点における結束度とする。ここでは、次のような結果が得られる。

【0105】

(5)

(6)

(7)

る場合には、「共通語彙比による結束度」と称することにする。

【0109】次に、図20は、ステップS43で記録された結束度の系列を示している。ここでは、窓幅の  $1/4$  が刻み幅  $t_{ic}$  として用いられており、文書領域  $a1 \sim a11$  は、刻み幅  $t_{ic}$  に対応する一定幅の領域である。また、 $c1$  は、文書中の  $a4$  と  $a5$  の境界を基準点として計算した、窓幅  $w$  の結束度を表す。すなわち、 $c1$  は、文書領域  $a1 \sim a4$  の部分を左窓の範囲とし、 $a5 \sim a8$  の部分を右窓の範囲として計算された結束度である。

【0110】次の  $c2$  は、窓を  $t_{ic}$  分だけ右へずらして計算された結束度を表し、 $a5$  と  $a6$  の境界を基準点とする窓幅  $w$  の結束度である。このようにして、窓を  $t_{ic}$  分ずつ順に右へずらして計算した  $c1, c2, c3, c4, \dots$  を、文書の冒頭から末尾へ向かう窓幅  $w$  の結束度系列と呼んでいる。

【0111】図21は、上述の単語認定結果において、文書の冒頭から各基準点までの間に出現した内容語の延べ数を横軸にとり、640語の窓幅の結束度系列をプロットしたグラフである。例えば、図20の結束度  $c2$  の場合は、 $a1 \sim a5$  の領域中の内容語の延べ数が、文書における基準点の位置となる。ここでは、640語の窓幅の  $1/8$  (80語) を刻み幅  $t_{ic}$  として、文書の冒頭から末尾に向かって結束度を計算している。

【0112】ここで、 $w_{min}$  を固定し、窓幅  $r$  として2の累乗の値を用いるのが計算効率上望ましいとした理由について説明する。窓幅  $r$  として2の累乗の値を用いるのが望ましいのは、次の理由による。各窓幅の結束

度の計算においては、文書中のそれぞれの位置で、その位置の前後に設定した 2 つの窓内の領域とそれらを合わせた領域の 3 種類の領域に出現する語彙を調べる必要がある。例えば、共通語彙比による結束度を用いる場合には、この 3 種類の領域に出現する語彙の異なり語数を集計する必要があるが、余弦測度による結束度を用いる場合には、この 3 種類の領域に出現する語彙の出現頻度を集計する必要がある。

【0113】図 19 では、左窓と右窓の各々の中の語彙数と、これらに共通する共通語彙の数を集計しているが、これらの 2 つの窓を合わせた領域中の語彙数は、左窓中の語彙数と右窓中の語彙数の和から共通語彙数を差し引いた結果に一致する。したがって、図 19 の集計は、上述した 3 種類の領域に出現する語彙数を集計する処理と同値であり、必要な計算量はほとんど変わらない。

【0114】このとき、 $r$  を 2 の累乗の値にしておくと、小さな窓幅の結束度の計算のために集計した語彙数（または出現頻度）を、大きい窓幅の結束度の計算でも利用できるようになる。例えば、 $r$  として 2 を用いると、窓幅  $w_1$  の結束度の計算において前後の窓を合わせた領域で集計した語彙数（または出現頻度）が、窓幅  $w_0$  の結束度の計算における片方の窓内における語彙数（または出現頻度）としても使用できることになる。

【0115】また、 $w_{min}$  を固定しておくことが望ましいのは次の理由による。 $w_{min}$  を固定し、窓幅比  $r$  として 2 の累乗の値を用い、さらに結束度計算の刻み幅  $t$  は各窓幅の  $1/n$  ( $n$  は整数) としておくと、要約粒度決定部 25 が  $Wd$  を数えるために文書全体を走査する際に、文書を  $w_{min}/n$  の領域に分割して、結束度系列の計算に便利な形に変換することができる。

【0116】例えば、各出現語彙を、ハッシュ表などを用いて、語彙の異なりを識別する語彙番号に変換（数値化）し、 $w_{min}/n$  の各領域に、出現語彙の語彙番号とその出現頻度を記録しておくことなどが可能になる。こうしておけば、少なくとも結束度系列の計算においては、原文書にアクセスする必要がなくなるので、計算効率性が向上する。

【0117】また、一般的な OS（オペレーティングシステム）は、原文書の中でアクセスしなくても、原文書の物理的な大きさ（バイト数）を簡単に取得できる機能を持っているのが普通である。

【0118】このような OS 上では、最初に、原文書の物理的な大きさで最大窓幅の大体の大きさ（例えば、上限）の見当をつけおき、 $Wd$  を数えるために文書全体を走査する際に、結束度系列の計算も同時に行うように工夫することも考えられる。この方法によれば、利用可能な 1 次メモリの容量が小さい環境でも、原文書へのアクセス回数を減らすことができる。その他にも、計算上の色々な工夫が考えられる。

【0119】次に、話題構成認定部 26 は、サブモジュールの話題境界候補区間認定部 27 を使って、それぞれの窓幅の結束度系列を解析し、結束度の低い区間を話題境界候補区間として認定する（ステップ S 44）。

【0120】図 21 に示したように、結束度系列における極小点は、実際の話題境界（点線で示した際の境界）に対応することが多いが、すべての極小点が話題境界に対応するわけではない。話題境界候補区間認定部 27 は、結束度系列の極小点を手掛かりに、それぞれの結束度系列の窓幅程度の大きさの話題のまとまりの境界位置を区間推定する。本実施形態では、この処理を、移動平均法を用いて実現している。

【0121】次に、話題構成認定部 26 は、異なる窓幅の結束度系列に基づいて求めた話題境界候補区間を統合し、大きな窓幅の結束度系列から得られた大きな話題に関する境界と、小さい窓幅の結束度系列からのみ得られる小さい話題に関する境界とを区別して出力する（ステップ S 45）。これにより、話題構成認定処理が終了する。

【0122】ここで、出力される最終的な話題境界は、統合された話題境界候補区間のうち最も小さい窓幅、すなわち最小窓幅の話題境界候補区間を使って認定される。最終的な話題境界の認定に最小窓幅の話題境界候補区間を使う理由は、大きな窓幅の結束度系列は、窓位置の移動に対して鈍感であり、それだけを用いて認定すると、境界位置を十分精密に求めることができないからである。

【0123】次に、図 17 のステップ S 44 における話題境界候補区間認定処理について、図 20 および図 2 を使って説明する。ここで用いられる移動平均法は、株価の変動などの統計的分析方法である時系列分析（time series analysis）において、細かい変動を取り除いて大局的な傾向を把握するために使われている。本実施形態では、結束度系列の移動平均値を細かい変動を無視するために用いるだけでなく、それを移動平均の開始点における右結束力および移動平均の終了点における左結束力とみなすことで、話題境界候補区間（低結束度の区間）認定のための直接的な手掛かりとしている。

【0124】図 20 は、前述したように、結束度の系列  $c_1 \sim c_4$  と文書領域  $a_1 \sim a_{11}$  との関係を示している。結束度系列の移動平均値とは、例えば、 $(c_1 + c_2) / 2$ （2 項の移動平均）、 $(c_1 + c_2 + c_3) / 3$ （3 項の移動平均）、 $(c_1 + c_2 + c_3 + c_4) / 4$ （4 項の移動平均）のように、結束度系列において連続する  $n$  個の値を算術平均した値である。

【0125】図 21 は、図 20 の結束度系列の移動平均の例と文書領域との関係を示している。ここでは、移動平均の例として、図 20 の結束度の 2 項～4 項の移動平均が示され、それぞれの移動平均に関わる結束度の計算において、各文書領域が使用された回数が示されている。

る。このうち、下線を付けた値は、対応する文書領域が移動平均に関わるすべての結束度の計算に用いられていることを表す。

【0126】例えば、左上角の値“1”は、 $c1 \sim c4$ までの4項の移動平均において、文書領域a1が一度だけ左窓の一部として扱われたことを示している。また、その右の値“2”は、 $c1 \sim c4$ までの4項の移動平均において、文書領域a2が2回左窓の一部として扱われたことを示している。他の使用回数についても、同様である。

【0127】結束度は境界の前後の部分の結び付きの強さを表す指標であるので、領域a1を左窓に含んで得られた結束度c1を用いて計算された移動平均値も、領域a1が右の方向に結び付いているかどうかを示す指標の1つと考えられる。

【0128】言い換えれば、移動平均値は、移動平均をとった結束度の左窓部分の領域( $c1 \sim c4$ の4項平均に対してはa1～a7)が右方向に引っ張られる強さ(右結束力)の指標になっていると言える。一方、逆に、移動平均をとった結束度の右窓部分の領域( $c1 \sim c4$ の4項平均に対してはa5～a11)が左方向に引っ張られる強さ(左結束力)の指標になっているとも言える。

【0129】ここで、結束力の指標とそれぞれの文書領域との関連性を考察すると、結束度の計算においてより多く窓に含まれていた領域との関連が強いと考えられる。また、語彙的結束度は、一般に、近傍で繰り返される語彙に基づくものほど強いと考えられるので、移動平均をとった結束度の基準点(左右の窓の境界位置)に近い位置にある領域ほど関連が強いとも言える。

【0130】例えば、図22の4項の移動平均については、結束度の基準点は、a4とa5の境界、a5とa6の境界、a6とa7の境界、およびa7とa8の境界の4つである。この場合、a4は最も多く左窓に含まれており、かつ、これらの基準点に最も近いことが分かる。また、a8は最も多く右窓に含まれており、かつ、これらの基準点に最も近いことが分かる。したがって、移動平均値と最も関連の強い領域は、左窓についてはa4、右窓についてはa8となる。

【0131】同様にして、3項の移動平均と最も関連の強い領域を選ぶと、左窓についてはa4、右窓についてはa7となり、2項の移動平均と最も関連の強い領域を選ぶと、左窓についてはa4、右窓についてはa6となる。これらの領域の使用回数は、図22では斜線を付けて示されている。

【0132】以上の考察に基づき、話題境界候補区間認定部27は、結束度の移動平均値を、移動平均をとった領域内の最初の基準点における右結束力および最後の基準点における左結束力の指標として取り扱う。例えば、 $c1 \sim c4$ の4項の移動平均値は、a4とa5の境界に

おける右結束力およびa7とa8の境界における左結束力となる。

【0133】図23は、話題境界候補区間認定部27による話題境界候補区間認定処理のフローチャートである。候補区間認定部27は、まず、話題構成認定部26から結束度系列の刻み幅ticを受け取り、ユーザから移動平均の項数nを受け取る(ステップS51)。

【0134】これらのパラメータの値の目安は、刻み幅ticについては、例えば、窓幅wの $1/8 \sim 1/10$ 程度の大きさであり、項数nについては、 $w/tic$ の半分(4～5)程度である。また、移動平均をとる領域の最初の基準点から最後の基準点までの隔たりを、 $(n-1) * tic$ により計算して、それを移動平均の幅d(語)とする。

【0135】次に、文書中の各位置pについて、 $p-p+d$ の範囲内で結束度の移動平均をとり、平均値を位置pにおける右結束力として記録する(ステップS52)。この値は、同時に、移動平均をとった範囲の終了位置p+dにおける左結束力としても記録される。

【0136】次に、記録された右結束力をもとに、文書中の冒頭から末尾に向かって各位置における右結束力と左結束力の差(右結束力-左結束力)を調べ、その値が負から正に変化する位置を負の結束力拮抗点として記録する(ステップS53)。

【0137】負の結束力拮抗点とは、その位置の左では左結束力が優勢であり、その位置の右では右結束力が優勢であるような点である。したがって、この点の左右の部分は意味的な結び付きが弱いと考えられ、負の結束力拮抗点は話題境界の候補位置となる。

【0138】次に、認定された結束力拮抗点の直前のd語以内の範囲で、右結束力が最小となる位置mpを求め、区間[mp, mp+d]を話題境界候補区間と認定して(ステップS53)、処理を終了する。

【0139】ここで、左右の結束力の差に基づいて話題境界候補区間を認定する意味を、図24を使って説明する。図24は、図21の5000語の近傍(4600語～5400語付近)における320語幅の窓における左結束度と左右の結束力の分布を示している。刻み幅ticとしては、窓幅の $1/8$ を採用している。

【0140】図24において、記号◇でプロットした折れ線グラフは、結束度Cの系列を表し、記号□でプロットした折れ線グラフは、右結束力FCの系列を表し、記号×でプロットした折れ線グラフは、左結束力BCの系列を表す。話題境界候補区間と結束力拮抗点を表す2重矩形で示された領域については、後述することにする。

【0141】また、点線で示されたb p1、b p2、b p3は、左右の結束力の差が0になる3つの点(結束力拮抗点)を表す。最初の点b p1の左側では、左結束力が右結束力より優勢であり、その右側から次の点b p2までは、右結束力が左結束力より優勢である。さらに、

その右側から最後の点  $b p 3$  までは、左結束力が右結束力より優勢であり、その右側では、右結束力が左結束力より優勢である。

【0142】したがって、 $b p 1$  と  $b p 3$  は、右結束力と左結束力の差が負から正に変化する負の結束力拮抗点であり、 $b p 2$  は、その差が正から負に変化する正の結束力拮抗点である。

【0143】このような結束力の変化から、最初の点  $b p 1$  の左側の領域は、それより左側のいずれかの部分と比較的強い結束性を示しており、真中の点  $b p 2$  の両側の領域は、 $b p 2$  に向かって強い結束性を示しており、最後の点  $b p 3$  の右側の領域は、それより右側のいずれかの部分と比較的強い結束性を示していることが分かる。実際、左右の結束力と共にプロットした結束度は、 $b p 1$  と  $b p 3$  の近傍で極小値をとり、 $b p 2$  の近傍で極大値をとっている。このように、左右の結束力の変化と結束度の変化は密接に関連している。

【0144】例えば、図 24 の結束力拮抗点  $b p 3$  の近傍の曲線で囲まれた部分 P 3 は、結束度が極小となる部

$$F C(p-d) = B C(p)$$

であり、拮抗点  $b p 3$  では左右の結束力が等しいので、

$$F C(b p 3-d) (= B C(b p 3)) = F C(b p 3) \quad (9)$$

が成り立つ。したがって、拮抗点  $b p 3$  の直前の点の右結束力が  $b p 3$  の値より小さければ、 $b p 3-d$  から  $b p 3$  までの範囲、すなわち、 $b p 3$  から左に  $d$  語以内の範囲に、右結束力の極小値が存在することになる。

$$F C(b p 3-d-1) = B C(b p 3-1)$$

$$> F C(b p 3-1)$$

$$\geq F C(b p 3)$$

(10)

が成り立つ。さらに、 $b p 3$  の右側において、

$$F C(b p 3) < F C(b p 3+1) \quad (11)$$

または、

$$F C(b p 3) \geq F C(b p 3+1) \quad (12)$$

が成り立つ。(11) 式が成り立つとき、(10)、(11) より、 $b p 3-d$  から  $b p 3$  までの範囲内に、

$$F C(b p 3-d+1) = B C(b p 3+1)$$

$$< F C(b p 3+1)$$

$$\leq F C(b p 3)$$

(13)

となる。したがって、(10)、(13) 式より、 $b p 3-d$  から  $b p 3$  までの範囲内に、右結束力の極小値が存在することになる。

【0148】図 25 は、図 17 のステップ S 45 において行われる話題境界認定処理のフローチャートである。話題構成認定部 26 は、まず、認定された話題境界候補区間を、認定に使った結束度系列の窓幅と、話題境界候補区間内の結束力拮抗点の文書における出現位置とによってソートしてまとめ、話題境界候補区間データの系列  $B(i)[p]$  を作成する (ステップ S 61)。

【0149】ここで、制御変数  $i$  は、窓幅  $w_i$  の結束度系列により認定されたことを表す系列番号であり、制御

分の 1 つである。このため、この部分 P 3 の移動平均 (ここでは、4 項平均) の値も、P 4 および P 5 における結束力が示しているように、通常は、極小値をとる。ただし、移動平均をとる領域より狭い範囲で細かい変動がある場合には、移動平均の平滑化作用により、移動平均値すなわち結束力が極小値をとらないこともある。

【0145】また、右結束力は移動平均値を移動平均をとる領域の開始位置に記録した指標であるので、右結束力の極小位置は結束度の極小位置の左になる。同様の理由により、左結束力の極小位置は結束度の極小位置の右になる。そして、結束度の変動が十分に大きければ、移動平均をとる領域内に結束力拮抗点が生成されることになる。

【0146】また、負の結束力拮抗点の直前の  $d$  語以内の範囲に右結束力の極小点が存在することは、次のようにして保証される。まず、ある点  $p$  における右結束力、左結束力を、それぞれ、 $F C(p)$ 、 $B C(p)$  とおくと、結束力の定義から、

(8)

【0147】また、拮抗点  $b p 3$  の直前の点の右結束力が  $b p 3$  の値より小さくない場合は、 $b p 3$  の左側において、

右結束力の極小値が存在することになる。また、(12) 式が成り立つとき、

変数  $p$  は、系列内の各話題境界候補区間を表すデータ番号である。実際には、 $i$  は、窓幅の大きい順に 0, 1, 2, ... のような値をとり、 $p$  は、結束力拮抗点の出現順に 1, 2, ... のような値をとる。それぞれのデータ  $B(i)[p]$  は、次のような要素データを含む。

【0150】 $B(i)[p].level$ : 話題境界のレベル。初期値は  $i$ 。

$B(i)[p].range$ : 話題境界候補区間。(開始位置、終了位置) の組。

【0151】 $B(i)[p].bp$ : 結束力拮抗点。(開始位置、終了位置) の組。

ここで、結束力拮抗点は理論的には点であるが、前述の

ように、右結束力と左結束力の差の符号が反転する地点を拮抗点として認定しているので、差が負の点を開始位置とし、差が正の点を終了位置とする小さな区間で表される。この区間の幅は、多くの場合、話題境界候補区間の認定に用いられた刻み幅  $t_{ic}$  に一致する。

【0152】  $B(i)[p]$ 、 $b_p$  としては、(開始位置、終了位置)の組を用いる代わりに、別法として、次

$$b_p = (DC(r_p) * l_p - DC(l_p) * r_p) / (DC(r_p) - DC(l_p)) \quad (14)$$

次に、話題構成認定部 26 は、出力対象とする話題境界のレベルの範囲を決定する(ステップ S 62)。出力対象とする話題境界が、基本窓幅  $w_1$ 、基本窓幅よりひとまわり大きい窓幅(最大窓幅)  $w_0$ 、および基本窓幅よりひとまわり小さい窓幅  $w_2$  の 3 種類の窓幅によって認定された話題境界候補区間に対応する場合は、 $L = \{0, 1, 2\}$  となる。

【0154】基本窓幅  $w_1$  による話題境界だけでなく、それに準ずる大きさの窓幅  $w_0$ 、 $w_2$  による話題境界も出力対象とするのは、次に行われる重要語抽出処理で話題に特徴的な語彙を選択する際に、これらの話題境界が使われるからである。窓幅比  $r$  が 2 で、基本窓幅  $w_1$  が 1280 語の場合、 $w_0 = 2560$  (語)、 $w_2 = 1280$  (語)、および  $w_2 = 640$  (語) の 3 種類の窓幅の話題境界が出力対象となる。

【0155】次に、話題構成認定部 26 は、窓幅の異なる話題境界候補区間データを統合する処理を行う。ここでは、1 つの系列に属する  $B(i)[p]$  をまとめて  $B(i)$  と記し、さらに、次のような表記法を用いて、以下の処理を説明する。

【0156】 $\cdot w_i : B(i)$  の系列番号  $i$  に対応する窓幅。

$\cdot d_i : B(i)$  に対応する話題境界候補区間の幅(移動平均の幅)。

$\cdot i_e$  : 最小窓幅  $w_{min}$  に対応する系列番号。

【0157】 $\cdot |B(i)| : B(i)$  におけるデータ番号  $p$  の最大値。

まず、処理対象を表す系列番号  $i$  を 0 に初期化する(ステップ S 63)。これにより、最大窓幅  $w_0$  による話題境界候補区間の系列が処理対象に設定される。次に、処理対象の系列  $B(i)$  に含まれるデータ  $B(i)[p]$  のうち、出力対象外のデータを取り除く(ステップ S 64)。すなわち、 $B(i)[p].level \in L$  となるデータ  $B(i)[p]$  だけを残し、その他のデータを  $B(i)$  から除外する。

【0158】そして、 $i$  をインクリメントしながら、 $i+1 \leq i_e$  である限り、 $B(i+1)$  を統合対象の系列とする統合処理を行う。この統合処理では、処理対象系列中のそれぞれの話題境界候補区間データ  $B(i)$

$[p] (p=1, \dots, |B(i)|)$  について、それと同じ付近を境界候補としている統合対象系列中のデ

のような位置データを用いてもよい。まず、結束力拮抗点の開始位置  $l_p$  と終了位置  $r_p$  における(右結束力-左結束力)の値を、それぞれ、 $DC(l_p)$  と  $DC(r_p)$  とする。そして、左右の結束力が 0 になる点  $b_p$  を、次式により補間して求めて、それを  $B(i)[p]$ 、 $b_p$  とする。

【0153】

ータ  $B(i+1)[q]$  が検出され、両者のデータが統合される。

【0159】この処理を途中で打ち切っても可能であるが、大きい窓幅に対応する系列で処理を打ち切ると境界位置の精密度が落ちることになる。また、この処理にはそれぞれの計算量は必要ないので、通常は、最小窓幅に対応する系列まで処理を繰り返す。

【0160】具体的な手順は以下の通りである。まず、 $i+1$  と  $i_e$  を比較し(ステップ S 65)、 $i+1 \leq i_e$  であれば、 $p$  に 1 を代入して(ステップ S 66)、 $p$  と  $|B(i)|$  を比較する(ステップ S 67)。 $p \leq |B(i)|$  であれば、図 26 の統合処理を行い(ステップ S 68)、 $p = p+1$  において(ステップ S 69)、ステップ S 67 以降の処理を繰り返す。そして、ステップ S 67 において、 $p$  が  $|B(i)|$  を越えれば、 $i = i+1$  において(ステップ S 70)、ステップ S 64 以降の処理を繰り返す。

【0161】そして、ステップ S 65 において、 $i+1$  が  $i_e$  を越えれば、統合処理を終了する。ここで、系列  $B(i_e)$  のそれぞれのデータ  $B(i_e)[p]$  について、 $B(i_e)[p].range$  の区間内で窓幅  $w_{ie}$  の結束度が最小となる位置  $mp$  を求め、 $mp$  と  $B(i_e)[p].level$  とを対応付けて出力する(ステップ S 71)。これにより、話題境界認定処理が終了する。

【0162】次に、図 26 の統合処理について説明する。話題構成認定部 26 は、まず、統合対象系列中のデータ  $B(i+1)[q] (q=1, \dots, |B(i+1)|)$  の中から、 $B(i+1)[q].b_p \cap B(i)[p].range \neq \emptyset$  であり、かつ、 $B(i+1)[q].b_p \cap B(i)[p].b_p$  が点で指定されている場合は、代わりに、 $B(i+1)[q].b_p \in B(i)[p].range$  という条件が用いられる。

【0163】ここで、 $B(i+1)[q].b_p \cap B(i)[p].range \neq \emptyset$  という条件は、 $B(i)[p]$  の話題境界候補区間と  $B(i+1)[q]$  の結束力拮抗点の区間とが、少なくとも部分的に重複していることを表す。 $B(i+1)[q].b_p$  が点で指定されている場合は、代わりに、 $B(i+1)[q].b_p \in B(i)[p].range$  という条件が用いられる。

【0164】図 27 は、統合対象データの選択例を示し

ている。図27において、記号◇でプロットした折れ線グラフは、処理対象に対応する640語幅の窓による右結束力の系列を表し、記号+でプロットした折れ線グラフは、640語幅の窓による左結束力の系列を表す。また、記号□でプロットした折れ線グラフは、統合対象に対応する320語幅の窓による右結束力の系列を表し、記号×でプロットした折れ線グラフは、320語幅の窓による左結束力の系列を表す。

【0165】また、2重矩形で示された領域のうち、大きな矩形領域が話題境界候補区間に対応し、それに含まれている小さな矩形領域が結束力拮抗点に対応する。ここでは、処理対象データB(i)[p]の話題境界候補区間と統合対象データのデータB(i+1)[q]を照合する際に、B(i)[p]の話題境界候補区間の幅を、前述の[m<sub>p</sub>, m<sub>p</sub>+d]よりt<sub>ic</sub>/2だけ左右に拡大し、[m<sub>p</sub>-t<sub>ic</sub>/2, m<sub>p</sub>+d+t<sub>ic</sub>/2]としている。t<sub>ic</sub>/2は、B(i+1)[q]に対応する結束度の刻み幅である。

【0166】これは、話題境界候補区間の認定精度が結束度系列の刻み幅t<sub>ic</sub>に依存するため、m<sub>p</sub>の本当の値は(m<sub>p</sub>-t<sub>ic</sub>, m<sub>p</sub>+t<sub>ic</sub>)の間と推定されるからである。したがって、処理対象データの話題境界候補区間を広めにとった場合には、(m<sub>p</sub>-t<sub>ic</sub>, m<sub>p</sub>+d+t<sub>ic</sub>)の範囲となる。

【0167】ここでは、統合対象データの結束度の刻み幅がt<sub>ic</sub>/2であるため、[m<sub>p</sub>-t<sub>ic</sub>/2, m<sub>p</sub>+d+t<sub>ic</sub>/2]を処理対象データの話題境界候補区間としている。このような話題境界候補区間を設定すれば、その幅はd+t<sub>ic</sub>=n\*t<sub>ic</sub>となる。一方、t<sub>ic</sub>=w/8、n=4であるから、話題境界候補区間の幅は、丁度、窓幅の半分w/2となる。

【0168】例えば、処理対象データをB(2)[6]とすると、その話題境界候補区間B(2)[6].rangeには、統合対象系列の2つのデータの結束力拮抗点B(3)[11].bpとB(3)[12].bpが含まれている。このため、B(3)[11]とB(3)[12]が統合対象データの候補となる。これらのうち、B(3)[12].bpの方が、処理対象データの結束力拮抗点B(2)[6].bpにより近いので、B(3)[12]が統合対象データとして選択される。

【0169】次に、話題構成認定部26は、統合対象データが選択できたかどうかを判定し(ステップS82)、統合対象データが選択できれば、ステップS84の処理を行う。ステップS81において、条件を満たすデータが見つからなかった場合には、処理対象データを認定するときに使った結束度を手掛かりに、擬似的な統合対象データを作成し、B(i+1)の系列に挿入する(ステップS83)。そして、ステップS84の処理を行う。

【0170】ステップS83では、まず、B(i)

[p].rangeの範囲内で、窓幅w<sub>i</sub>の結束度が最小となる位置m<sub>p</sub>を求める。次に、B(i+1)[q].bp=[m<sub>p</sub>, m<sub>p</sub>]、B(i+1)[q].range=[m<sub>p</sub>-d<sub>i</sub>/2, m<sub>p</sub>+d<sub>i</sub>/2]と設定して、m<sub>p</sub>に対応する新たなデータB(i+1)[q]を作成する。

【0171】そして、系列B(i+1)の中で、B(i+1)[q-1].bp<m<sub>p</sub>かつB(i+1)[q+1].bp>m<sub>p</sub>となるような位置に、作成したデータB(i+1)[q]を挿入する。これにより、疑似的な統合対象データのデータ番号qが決定され、それ以降の既存データのデータ番号は書き換えられる。ここで、擬似的な話題境界候補区間データを作成するのは、以降の処理において統合探索範囲を狭め、精密な境界認定を行うためである。

【0172】例えば、図27のB(2)[6]を処理対象データとすると、通常の統合対象データの話題境界候補区間B(3)[12].rangeの幅は、d<sub>3</sub>(160語)である。このときに、もし、B(3)[11]とB(3)[12]のいずれも存在しなかった場合には、図28に示すように、B(2)[6].rangeの範囲内における窓幅w<sub>i</sub>(640語)の結束度が最小値をとる位置m<sub>p</sub>を求める。

【0173】そして、その近傍にB(3)[10].rangeなどの通常の話題境界候補区間と同じ幅d<sub>3</sub>のB(3)[q].rangeを持つ疑似的なデータB(3)[q]を作成する。これにより、ステップS84の処理において、B(2)[6].rangeの幅d<sub>2</sub>(320語)をd<sub>3</sub>(160語)に絞り込むことができる。

【0174】この操作は、処理対象データの話題境界候補区間において、結束度の極小点が明確に1点に決まる場合には、大抵の場合有効である。しかし、話題境界候補区間において結束度ほとんど変動が見られない場合には、その話題境界候補区間を縮小せずに、そのまま用いた方がよいこともある。ただし、経験的には、話題境界候補区間において結束度がほとんど変動しないような状況は、あまり多く現れない。

【0175】ステップS84では、統合対象データの話題境界レベルB(i+1)[q].levelを処理対象データの話題境界レベルB(i)[p].levelに変更して、処理対象データB(i)[p]と統合対象データB(i+1)[q]の情報を統合する。この処理は、統合対象データB(i+1)[q]の話題境界レベルを小さくすることに対応する。例えば、図27の統合対象データB(3)[12]の場合は、B(3)[12].level=B(2)[6].levelとなる。

【0176】これにより、次にステップS64の処理を行うとき、新たに処理対象となる系列B(i+1)の中



のデータのうち、少なくとも統合対象データ B (1+1) [q] は除外されずに残されることになる。したがって、処理対象データを統合対象データに順次置き換えながら、話題境界候補区間を徐々に絞り込んでいくことができる。

【0177】最終的には、統合対象データが系列 B (1e) から選択され、それぞれの統合対象データ B (1e) [p] が、ステップ S71 の処理が行われる。こうして出力された位置 mp が、その話題境界レベル B (1e) [p]、level における話題境界として認定される。

【0178】図 29 は、こうして得られた話題境界の認定結果を示している。図 29 において、2560 語、1280 語、640 語の各窓幅に対応して 2 重矩形で示された領域のうち、大きな矩形領域が話題境界候補区間に対応し、それに含まれている小さな矩形領域が結束度抵抗点に対応する。B (0)、B (1)、B (2) は、それぞれ、2560 語、1280 語、640 語の各窓幅に対応する系列を表し、2 重矩形に添えられた番号

[1], [2], ... などは、各系列内のデータ番号を表す。

【0179】また、上にある矩形領域ほど大きな窓幅（小さな話題境界レベル）に対応し、下にある矩形領域ほど小さな窓幅（大きな話題境界レベル）に対応する。そして、記号 \* の付いた棒グラフは、最終的に求められた話題境界の位置を表す。

【0180】後述する重要語抽出処理では、大きな窓幅の結束度に基づいて認定された境界ほど（棒グラフが長いほど）、大きな話題のまとまりに関する境界（話題境界レベルの小さな境界）であるとみなされる。また、小さな窓幅の結束度に基づいて認定された境界ほど（棒グラフが短いほど）、小さな話題のまとまりに関する境界（話題境界レベルの大きな境界）であるとみなされる。

$$\text{再現率} = (\text{正解数} / \text{節境界数}) * 100 (\%)$$

$$\text{適合率} = (\text{正解数} / \text{認定境界数}) * 100 (\%)$$

ここで、節境界数は、各窓幅における正解データの数を表し、認定境界数は、各窓幅の話題境界レベルに対応する認定境界の数を表し、正解数は、各窓幅において、正解データとの隔たりが 4 語以内であるような認定境界の数を表す。

【0187】例えば、4. 4 節の先頭の境界は、4. 3 節 (6, 067 語) と 4. 4 節 (6, 670 語) の間にあり、小さい方の節の大きさは 6, 067 語である。これは最大窓幅の 2, 560 語より大きいため、4. 4 節の先頭の境界は、すべての窓幅において正解データとして扱われる。

【0188】また、4. 4. 1 節の先頭の境界は、4. 4 節の先頭から 4. 1 節の先頭までの部分 (115 語) と 4. 4. 1 節 (2, 643 語) の間にあり、小さい方の節の大きさは 115 語である。したがって、小

【0181】図 29 の認定結果では、4. 3 節の開始位置に対応する境界 P11 より、その前の境界 P12

(4. 2. 2 (3) 節の開始位置に対応) の方が、大きな話題の境界であると認定されている。このような若干の食い違いはあるものの、大旨、大きな窓幅によって認定された境界ほど大きな話題の切れ目に対応するという傾向にあることが見てとれる。

【0182】また、図 30 は、共通語彙比による結束度の代わりに、余弦測度による結束度を用いた場合の話題境界の認定結果を示している。図 30 においても、大旨、図 29 と同様の傾向が見てとれる。

【0183】図 31 から図 36 までは、各窓幅の結束度を手掛かりに認定された話題境界（認定境界）の特徴を表すデータの集計結果を示している。このうち、図 31 から図 36 までは、(7) 式により求めた共通語彙比による結束度を使った場合の結果を表し、図 34 から図 36 までは、余弦測度による結束度を使った場合の結果を表す。

【0184】図 31 と図 34 は、本発明の狙い通りに、窓幅に応じた大きさの話題のまとまりが認定できているかどうかを調べるために、認定境界の間隔を集計した結果である。これらの集計結果から、窓幅の 1~2 倍程度の間隔で話題境界が認定されていることが分かる。

【0185】また、図 32、33、35、36 は、窓幅程度の間隔で認定された境界が、実際に、その大きさ程度の話題のまとまりと対応しているかどうかを調べた結果である。図 32 と図 35 では、上述の要約対象文書に含まれる節の各境界について、その前後の節の大きさを調べ、小さい方の節の大きさが窓幅以上であるような境界を正解データとして、各窓幅毎に再現率と適合率を集計している。再現率と適合率は、次式により計算された。

$$\text{【0186】}$$

$$(15)$$

$$(16)$$

4. 1 節の先頭の境界は、80 語と 40 語の窓幅においてのみ、正解データとして扱われる。

【0189】また、図 33 と図 36 では、要約対象文書に含まれる節の各境界について、その前後の節の大きさを調べ、小さい方の節の大きさが窓幅の 1/2 以上であるような境界を正解データとして、(15)、(16) 式により再現率と適合率を集計している。

【0190】これらの結果を比較すると、共通語彙比による結果より、余弦測度による結果の方が若干精度が高い。一方、同じ窓幅の結束度に関しては、共通語彙比による結果の方が、認定境界が多めになっている。これは、共通語彙比による結束度の方が、余弦測度による結束度より、繰り返される語彙数の変化に敏感であることによるものと考えられる。

【0191】このため、共通語彙比による結束度、小

さい窓幅では局所的な特異点の影響を受けやすく、  
±4語（合わせて1文程度の大きさ）の精度においては、若干見劣りのする結果を与えている。逆に、大きな窓幅においては、共通語彙比による結束度は、余弦測度による結束度では感知できない変動を拾うことができたものと考えられる。

【0192】本発明の実施に当たっては、これらの性質と結束度の計算のためのコストとを考慮し、適切な結束度の計算方法を選択あるいは併用することが望ましい。一般に、共通語彙比による結束度は計算コストが比較的  
10 低いため、計算効率を重視する場合にはこれを用いることが推奨される。

【0193】次に、結束度と文書中の書式を用いて、話題境界をより精度高く認定する方法を説明する。話題境界候補区間は、図29に見られるように、実際の節境界を含んでいる確率が高い。このため、図37に示するような簡単な書式の特徴を手掛かりとして、認定境界の位置を微調整することで、認定結果の精度を上げることが可能である。

【0194】図37には、この調整で用いる書式ボタンと境界レベルの関係が示されている。書式ボタンとしては、節境界を認定する手掛かりとなる特徴的な文字列が、一般的なOSで用いられている正規表現（regular expression）記法により示されている。例えば、「外1」は、「4.1」などのようにピリオドで区切られ  
20 【0195】

【外1】  

$$\backslash d+ \backslash . \backslash d+ \backslash . [ \ ] + \$$$

【0196】た2つの数字で始まり、句点「。」を含む  
30 ない行を表す。また、境界レベルとしては、上述の話題境界レベルと同様に、小さいほど大きな話題境界に対応するような番号が割り振られている。例えば、4.1節などはレベル1の境界であり、空行（「外2」）はレベル4の境界となる。

【0197】

【外2】  

$$\backslash \$$$

【0198】図38は、このような特定の書式ボタンを用いた統合処理のフローチャートである。この統合処理  
40 は、図25のステップS68で行われる。図37のような書式ボタンと境界レベルの関係は、あらかじめユーザにより指定されるものとする。

【0199】話題構成認定部26は、まず、与えられた書式ボタンを参照しながら、統合対象データの話題境界候補区間B(i)[p]、range内を走査し、最も境界レベルが小さく、B(i)[p]、bpに最も近い節境界の位置hpを求める（ステップS91）。そして、統合対象系中のデータB(i+1)[q]（q=1, . . . , |B(i+1)|）の中から、hp<B  
50

(i+1)[q]、rangeとなるようなデータB(i+1)[q]を、統合対象データとして選択する（ステップS92）。

【0200】次に、統合対象データが選択できたかどうかを判定し（ステップS93）、統合対象データが選択できれば、ステップS95の処理を行う。ステップS92において、条件を満たすデータが見つからなかった場合には、節境界hpを用いて疑似的な統合対象データを作成し、B(i+1)の系列に挿入する（ステップS94）。そして、ステップS95の処理を行う。

【0201】ステップS94では、B(i+1)[q]、bp=[hp, hp]、B(i+1)[q]、range=[hp-dw/2, hp+dw/2]と設定して、hpに対応する新たなデータB(i+1)[q]を作成する。

【0202】そして、系列B(i+1)の中で、B(i+1)[q-1]、bp<hpかつB(i+1)[q+1]、bp>hpとなるような位置に、作成したデータB(i+1)[q]を挿入する。これにより、疑似的な統合対象データのデータ番号qが決定され、それ以降の既存データのデータ番号は書き換えられる。

【0203】ステップS95では、図26のステップS84と同様に、統合対象データの話題境界レベルB(i+1)[q]、levelを処理対象データの話題境界レベルB(i)[p]、levelに変更して、処理対象データB(i)[p]と統合対象データB(i+1)[q]の情報を統合する。

【0204】このような統合処理を採用した場合は、図25のステップS71において、結束度の最小位置mpを求める代わりに、図38のステップS91と同様に、B(i)[p]、range内で境界レベルの最も小さい節境界の位置hpを求める。そして、hpとB(i)[p]、levelを対応付けて出力する。

【0205】図38の統合処理によれば、要約対象文書に含まれる実際の書式ボタンを手掛かりとして話題境界が認定されるため、図26の統合処理に比べて、認定結果の精度が向上する。

【0206】次に、重要箇所特定部28の処理について説明する。重要箇所特定部28は、話題構成認定部26が認定した話題境界で区切られた3つのレベルの話題区間のうち、結束度の低いものを以降の要約処理の対象から除外する。

【0207】ここで、3つのレベルの話題区間とは、最大窓幅w0の結束度によって求められた話題境界により区切られた区間と、基本窓幅w1以上の窓幅の結束度によって求められた話題境界で区切られた区間と、基本窓幅の次に小さい窓幅w2（=w1/r）以上の窓幅の結束度によって求められた話題境界で区切られた区間の3種類の話題区間を指す。低結束度の区間を処理対象から除外するのは、このような区間は、例えば、項目を羅列

しただけの部分のように、内容が薄い部分であることが多いためである。

【0208】ここでは、ある話題区間の結束度を、話題の階層的構成における親の話題区間における結束度の平均値と比較することで、その話題区間が低結束度区間であるかどうかを判定する。具体的には、判定対象の話題区間を $b$ とし、 $b$ の窓幅を $w_b$ とし、話題区間 $b$ の中心

$$c < mc + a$$

ここで、 $a$ は、低結束度判定の感度を変更するためのパラメータであり、この値が大きいくほど、低結束度区間と判定される区間が増える。 $a$ としては、0または親話題区間 $a$ における $w_a$ の標準偏差などを用いることが望ましい。

【0210】図39および図40は、重要箇所特定部28による重要箇所特定処理のフローチャートである。重要箇所特定部28は、まず、文書全体を親話題区間として設定し（図39、ステップS101）、最大窓幅 $w_0$ により決められた話題区間から、低結束度区間を除外する（ステップS102）。ここで、文書全体を親話題区間とするのは、 $w_0$ に基づく話題境界より上位の話題境界は存在しないためである。

【0211】最大窓幅 $w_0$ の話題区間は、直接には以後の要約処理の対象とはならないが、 $w_0$ の話題区間が除外されると、その中に含まれる基本窓幅 $w_1$ の話題区間もすべて除外されるため、要約処理の対象となる話題区間が減少する。

【0212】次に、基本窓幅 $w_1$ に基づく話題区間から、低結束度区間を除外する。基本窓幅 $w_1$ の話題区間の親話題区間は、最大窓幅 $w_0$ の話題区間であるので、ステップS102の処理によって除外されなかった $w_0$ の話題区間を1ずつ取り出し、その中に含まれる $w_1$ の話題区間を除外する。

【0213】ここでは、まず、最大窓幅 $w_0$ の最初の話題区間を取り出し、それを親話題区間として（ステップS103）、基本窓幅 $w_1$ の話題区間から、低結束度区間を除外する（ステップS104）。次に、最大窓幅 $w_0$ の次の話題区間を取り出し、それを親話題区間とする（ステップS105）。そして、親話題区間が取り出せたかどうかを判定し（ステップS106）、それが取り出せた場合は、ステップS104以降の処理を繰り返す。

$$c0 = mc + a$$

次に、親話題区間 $a$ において、窓幅 $w$ の最初の話題区間を取り出し、それを処理対象話題区間とする（ステップS113）。そして、処理対象話題区間の中心付近における最大結束度 $c$ を求める（ステップS114）。

【0219】次に、 $c$ と $c0$ を比較し（ステップS115）、 $c < c0$ であれば、処理対象話題区間を要約処理の対象から除外する（ステップS116）。そして、親話題区間 $a$ において、窓幅 $w$ の次の話題区間を取り出

付近における窓幅 $w$ 。による結束度の最大値を $c$ とし、窓幅 $w_{b-1}$ の話題区間のうち $b$ を含むものを親話題区間 $a$ とし、 $a$ における窓幅 $w_a$ による結束度の平均値を $m$ と $c$ とする。そして、次のような関係が成り立てば、話題区間 $b$ を低結束度区間であると判定する。

【0209】

(17)

【0214】親話題区間が取り出せなかった場合は、基本窓幅 $w_1$ の話題区間の除外処理が終了したものとみなし、次に、基本窓幅の次に小さい窓幅 $w_2$ の話題区間の除外処理を行う。窓幅 $w_2$ の話題区間の親話題区間は窓幅 $w_1$ の話題区間であるので、ステップS104の処理によって除外されなかった $w_1$ の話題区間を1ずつ取り出し、その中に含まれる $w_2$ の話題区間を除外する。

【0215】 $w_2$ の話題区間を除外することは、要約処理の対象である窓幅 $w_1$ の話題区間の中から内容的にまとまりの薄い部分を取り除くことに対応する。これにより、基本窓幅 $w_1$ の話題区間の要約文として、余分な内容が抜粋されることを防止できる。

【0216】ここでは、まず、基本窓幅 $w_1$ の最初の話題区間を取り出し、それを親話題区間として（図40、ステップS107）、窓幅 $w_2$ の話題区間から、低結束度区間を除外する（ステップS108）。次に、基本窓幅 $w_1$ の次の話題区間を取り出し、それを親話題区間とする（ステップS109）。そして、親話題区間が取り出せたかどうかを判定し（ステップS110）、それが取り出せた場合は、ステップS108以降の処理を繰り返す。

【0217】親話題区間が取り出せなかった場合は、窓幅 $w_2$ の話題区間の除外処理が終了したものとみなし、処理を終了する。図41は、図39のステップS102、S104、および図40のステップS108において呼び出される話題区間除外処理のフローチャートである。話題区間除外処理のサブモジュールは、まず、話題区間の窓幅 $w$ とその親話題区間 $a$ を呼び出し元から受け取る（ステップS111）。そして、親話題区間 $a$ において、処理対象の窓幅 $w$ の結束度の平均値 $m$ を求め、次式により、判定の基準となる基準結束度 $c0$ を決定する（ステップS112）。

【0218】

(18)

し、それを処理対象話題区間とする（ステップS117）。 $c \geq c0$ であれば、処理対象話題区間を残したまま、ステップS117の処理を行う。

【0220】次に、処理対象話題区間が取り出せたかどうかを判定し（ステップS118）、それが取り出せた場合は、ステップS114以降の処理を繰り返す。そして、処理対象話題区間が取り出せなくなれば、処理を終了する。

【0221】図42は、図41のステップS114において呼び出される最大結束度計算処理のフローチャートである。最大結束度計算処理のサブモジュールは、まず、処理対象話題区間bとその話題区間の窓幅wを呼び出し元から受け取り（ステップS121）、話題区間bの大きさwと比較する（ステップS122）。

【0222】話題区間bの大きさがwより大きければ、話題区間bから、その両端w/2の部分を除外した区間における最大結束度を求め、その値をcとして記録して（ステップS123）、処理を終了する。また、話題区間bの大きさがw以下であれば、話題区間bの中心位置における結束度をcとして記録し（ステップS124）、処理を終了する。

【0223】図43は、 $\alpha=0$ として、重要箇所特定処理を上述の要約対象文書に適用した結果を示している。図43において、斜線部分P21、P22、およびP23は、窓幅 $w_1$ （1280語）の低結束度区間の除外処理により除外された話題区間を表す。また、横線は、窓幅 $w_0$ の各話題区間における窓幅 $w_1$ の結束度の平均値m cを表し、矢印は、窓幅 $w_1$ の各話題区間の中心付近において、最大結束度cに対応する点を表す。

【0224】例えば、4000語付近の斜線部分P21を見ると、矢印が指す極大値cは、明らかに平均値m cより低い値を示しているのが分かる。このため、この話題区間は要約対象から除外されている。他の斜線部分P

22、P23についても同様である。

【0225】また、ハッチングされた部分P24およびP25は、窓幅 $w_2$ （640語）の低結束度区間の除外処理により除外された話題区間を表す。この処理により除外されなかった部分、すなわち、P21、P22、P23、P24、およびP25以外の部分は、要約処理の対象となる重要箇所であると認定される。

【0226】図39および図40の重要箇所特定処理では、結束度が閾値より低い話題区間を除外することで重要な話題区間を特定するが、その代わりに、結束度が閾値以上の話題区間を抽出する処理を行っても、同様の結果が得られる。

【0227】次に、重要語抽出部29の処理について説明する。重要語抽出部29は、話題構成認定部26が認定し、重要箇所特定部28が絞り込んだ基本窓幅 $w_1$ および最大窓幅 $w_0$ の話題区間のそれぞれに特徴的に出現している内容語を選択し、話題区間との対応を付けてそれらを抽出する。

【0228】ここでは、ある内容語tの話題区間bにおける出現頻度（出現度数）が期待値を上回り、かつ、次式の数値尤度比Lが与えられた閾値（統計的有意水準に対応する $\chi$ 値）以上であるとき、内容語tは話題区間bに特徴的であると判定される。

【0229】

【数2】

$$L=2(F_{bt}\log \frac{F_{bt}}{E(F_{bt})} + (F_{at}-F_{bt})\log \frac{F_{at}-F_{bt}}{F_{at}-E(F_{bt})})$$

(19)

【0230】(19)式において、 $F_{bt}$ は、話題区間bにおける単語tの出現頻度を表し、 $F_{at}$ は、話題区間bの親話題区間aにおける単語tの出現頻度を表し、 $E(F_{bt})$ は、話題区間bにおける単語tの出現頻度の期待値を表す。 $E(F_{bt})$ は、親話題区間aにおける単語tの出現密度（出現確率）に、話題区間bの大きさを乗じて得られる。ここで、ある区間における単語の出現密度とは、単語の出現頻度と区間の大きさの比を意味する。

【0231】(19)式のLは、単語tの出現確率が話題区間bとそれ以外の領域との区別に対して独立であるかどうかに関する尤度比検定の値であり、この値が大きいかほど、単語の出現確率がその区別に依存していることを表す。Lの自由度 $\nu$ は1であるので、有意水準が10%なら閾値を6.63490とし、有意水準が5%なら閾値を7.87994とし、有意水準が1%なら閾値を10.8276とすればよい。あるいは、閾値を用いる代わりに、Lの大きい順に上位のいくつかの単語を重要語として抽出しても構わない。

【0232】なお、最大窓幅 $w_0$ の話題区間と基本窓幅

$w_1$ の話題区間が一致している場合や、窓幅 $w_0$ の話題区間のほとんどを1つの基本窓幅 $w_1$ の話題区間が占めている場合には、このような検定方法は必ずしもうまく機能しない。このため、bの直接の上位の話題区間（例えば、bを含む窓幅 $w_0$ の話題区間）の大きさがbの大きさの2倍未満の場合には、親話題区間として文書全体を用いることにする。

【0233】図44および図45は、重要語抽出部29による重要語抽出処理のフローチャートである。重要語抽出部29は、まず、パラメータ $\chi$ 値の統計的有意水準に対応する閾値hを、ユーザから受け取る（図44、ステップS131）。そして、文書全体を親話題区間候補a0とし、a0の大きさとa0に出現しているそれぞれの内容語wの出現頻度を求め、それぞれ、S0とFw0として記録する（ステップS132）。

【0234】次に、最大窓幅 $w_0$ の先頭の話題区間を取り出し、それを親話題区間候補a1とする（ステップS133）。そして、親話題区間候補a1の大きさとa1に出現しているそれぞれの内容語wの出現頻度を求め、それぞれ、S1とFw1として記録する（ステップS13

4)。

【0235】次に、 $F_w$ に記録されたそれぞれの内容語の出現頻度の対数尤度比を、 $a_0$ を基準として求め、それを閾値 $h$ と比較して重要語を抽出する(ステップS135)。そして、 $a_1$ において、基本窓幅 $w_1$ の最初の話題区間を取り出し、それを重要語抽出対象区間 $b$ とする(ステップS136)。

【0236】次に、重要語抽出対象区間 $b$ の大きさ、 $b$ に出現しているそれぞれの内容語 $w$ の出現頻度を求め、それぞれ、 $S_b$ と $F_w$ として記録する(図45、ステップS137)。そして、 $S_1$ と $S_b$ を比較する(ステップS138)。ここで、 $S_1 < S_b$ であれば、観話区間として $a_0$ を選択し(ステップS139)、 $S_1 \geq S_b$ であれば、観話区間として $a_1$ を選択して(ステップS140)、ステップS141の処理を行う。

【0237】ステップS141では、重要語抽出部29は、 $F_w$ として記録された各内容語の出現頻度の対数尤度比を求め、それを閾値 $h$ と比較して重要語を抽出する。そして、 $a_1$ において、基本窓幅 $w_1$ の次の話題区間を取り出し、それを重要語抽出対象区間 $b$ とする(ステップS142)。

【0238】次に、 $b$ が取り出せたかどうかを判定し(ステップS143)、それが取り出せた場合は、ステップS137以降の処理を繰り返す。そして、 $b$ が取り

$$E(F_{bt}) = Fat * S_b / Sa$$

ここで、 $F_{bt}$ が $E(F_{bt})$ より大きければ、(19)式により $t$ の尤度比 $L$ を求め(ステップS157)、それを閾値 $h$ と比較する(ステップS158)。 $L$ が $h$ 以上であれば、 $t$ を重要語として抽出する(ステップS159)。そして、 $F_w$ のリストから次の単語を取り出し、それを検定単語 $t$ として(ステップS160)、図46のステップS153以降の処理を繰り返す。

【0243】ステップS158において $L$ が $h$ 以上である場合は、話題区間 $b$ における単語 $t$ の出現頻度が、観話区間 $a$ における出現頻度にくらべて特異に大きいものとみなされ、 $t$ が重要語として抽出される。

【0244】ステップS154、S156、およびS158で判定結果がNOの場合は、 $t$ を重要語として抽出せずに、ステップS160以降の処理を行う。そして、ステップS153において $t$ が取り出せなかった場合は、すべての単語の検定が終了したもののみなし、処理を終了する。

【0245】図48は、上述の要約対象文書の見出しのうち、先頭的重要語抽出対象区間(窓幅 $w_1 = 1280$ 語の話題区間)に含まれる見出しの例を示しており、図49は、その区間から抽出された重要語を示している。ここでは、有意水準5%に対応する閾値を用いた。

【0246】次に、重要文選択部30の処理について説明する。重要文選択部30は、出題人による先頭の特

出せなくなると、次に、最大窓幅 $w_0$ の次の話題区間を取り出し、それを観話区間候補 $a_1$ とする(ステップS144)。

【0239】次に、 $a_1$ が取り出せたかどうかを判定し(ステップS145)、それが取り出せた場合は、図44のステップS134以降の処理を繰り返す。そして、 $a_1$ が取り出せなくなると、処理を終了する。

【0240】図46および図47は、図44のステップS135および図45のステップS141において呼び出される尤度比検定処理のフローチャートである。尤度比検定処理のサブモジュールは、まず、閾値 $h$ 、観話区間の大きさ $S_a$ ( $S_0$ または $S_1$ )、観話区間における単語の出現頻度 $F_w$ ( $F_w$ または $F_w$ )のリスト、検定対象話題区間の大きさ $S_b$ 、および検定対象話題区間における単語の出現頻度 $F_w$ のリストを呼び出し元から受け取る(図46、ステップS151)。

【0241】次に、 $F_w$ のリストから最初の単語を取り出し、それを検定単語 $t$ とする(ステップS152)。そして、 $t$ が取り出せたかどうかを判定し(ステップS153)、それが取り出せた場合は、 $F_{bt}$ を $L$ と比較する(図47、ステップS154)。

【0242】 $F_{bt}$ が1より大きければ、 $F_{bt}$ の期待値(理論値) $E(F_{bt})$ を次式により計算して(ステップS155)、それを $F_{bt}$ と比較する(ステップS156)。

$$(20)$$

平9-006777「文書要約装置およびその方法」に示された技術を応用して、要約の一部となる重要文を抽出する。

【0247】本実施形態における重要文選択処理の特徴は、要約を作成する単位である基本窓幅 $w_1$ の話題区間と、話語構成においてそれらの直接の上位に位置する最大窓幅 $w_0$ の話題区間の両方に対して、重要語が与えられるという点である。このように、階層的に構成された話題区間のそれぞれに重要語を与えて、異なる階層の重要語を併用して重要文を選択するという点において、重要文選択処理は先頭の方法とは異なっている。

【0248】基本窓幅 $w_1$ の話題区間に与えられる重要語は、その語に関連の深い文をその話題区間からのみ抜粋するために用いられる局所的な(ローカルな)重要語である。一方、最大窓幅 $w_0$ の話題区間に与えられる重要語は、その下位に位置する複数の要約対象話題区間のそれぞれから関連の深い文を抜粋するために用いられる大域的な重要語である。

【0249】特開平9-006777は、少ない抜粋量でも要約文書に重要語を幅広く取り入れることができる方法を示している。この方法によれば、多くの種類の重要語を含む要約文書を生成することができる。これは、1つの文を選択する度に、重要語のリストから、選択された文に含まれる語を取り除いているためである。

【0250】このリストは、ユーザの質問文の単語も含んでいるため、注目語リストと呼ばれている。ここで、注目語とは、文書の作成者が書こうとした内容を示すキーワード（端的には、見出しや強調語）と、ユーザが閲覧したいと思っている文書の事柄を示すキーワード（端的には、文書検索時にユーザが入力する質問文）の両方を含んでいる。

【0251】本実施形態においても、数十ページの文書を1ページ程度に要約することを念頭においているため、重要文を選択する度に、注目語リストを更新するとい

10 いう方法を踏襲する。  
【0252】図50および図51は、重要文選択部30による重要文選択処理のフローチャートである。重要文選択部30は、まず、最大窓幅 $w_0$ の先頭の話題区間を親話題区間 $a$ として取り出し（図50、ステップS161）、 $a$ に対応する重要語を親区間の注目語リスト $kw1$ に登録する（ステップS162）。そして、 $a$ の先頭部分に見出しが存在するかどうかを判定する（ステップS163）。見出しの判定には、例えば、図37に示した書式ボタンが用いられる。

【0253】見出しが存在する場合には、その見出しに印を付けて必須出力文（必ず要約に含める文）とし、見出しに含まれる内容語を抽出して注目語リスト $kw1$ に加える（ステップS164）。これにより、見出しに関連する文も、自動的に要約に採録されるようになる。

【0254】次に、 $a$ に含まれる要約対象の各話題区間を $b$ とし、 $b$ の注目語リスト $kw1n$ を作成する。注目語リスト $kw1n$ には、まず、各話題区間 $b$ に固有の重要語が登録され（ステップS165）、次に、親話題区間 $a$ の注目語リスト $kw1$ の注目語がマージさ

30 される（ステップS166）。ステップS163において見出し語が存在しない場合には、そのままステップS165以降の処理を行う。  
【0255】次に、 $a$ を親話題区間とするすべての $b$ を一度に処理して、それぞれの区間から要約に出力する文を1文ずつ選択する（図51、ステップS167）。ここで、同じ話題区間 $a$ を親に持つ子話題区間 $b$ を一度にまとめて処理するのは、 $a$ の注目語に関連する文なるべく多くの $b$ から採録することを意図しているためである。このとき、選択された文には、選択済であることを示す印が付けられる。

【0256】次に、文が選択できなかった話題区間 $b$ に対応する $kw1n$ を削除して、その区間に対する選択処理を終了する（ステップS168）。また、文が選択された話題区間 $b$ については、選択された文に含まれる注目語を、対応する注目語リスト $kw1n$ から削除する（ステップS169）。さらに、親話題区間 $a$ の注目語リスト $kw1$ のみから由来する注目語であって、話題区間 $b$ に固有の注目語ではないものが、別の話題区間 $b_x$ で選択された文に含まれている場合には、その注目

語を $b$ の注目語リスト $kw1n$ から削除する（ステップS170）。

【0257】次に、注目語リスト $kw1n$ が残っているかどうか、すなわち、まだ文を選択する余地のある話題区間 $b$ があるかどうかを判定する（ステップS171）。そのような注目語リストが残っている場合には、ステップS167以降の処理を繰り返す。このとき、ステップS169、S170の処理により空になってしまった注目語リスト $kw1n$ については、 $b$ に固有の注目語リストと $a$ の注目語リストをマージして、注目語リスト $kw1n$ の初期状態を復元しておく（ステップS172）。

【0258】また、ステップS171において、注目語リスト $kw1n$ が残っていない場合は、最大窓幅 $w_0$ の次の話題区間を親話題区間 $a$ として取り出す（ステップS173）。そして、親話題区間 $a$ を取り出せたかどうかを判定し（ステップS174）、それが取り出せた場合は、図50のステップS162以降の処理を繰り返す。

20 【0259】そして、親話題区間が取り出せなかった場合は、ステップS164で印を付けられた必須出力文とステップS167で選択された文とをマージし、出現順に並べて要約を作成して（ステップS175）、処理を終了する。作成された要約に、選択されなかった文の存在を示す印や段落境界などを挿入することで、要約の可読性を高めることも可能である。

【0260】ところで、ステップS167において文を選択できない場合としては、採録量の制約により文選択が打ち切られた場合、その時点の注目語リストに含まれている注目語（重要語）を含む文が見つからなかった場合などがある。後者の場合には、もう一度注目語リストの初期状態を復元して文選択を試みることで、選択可能な文の範囲を広げることができる。

【0261】また、ここでは、ユーザが入力した質問文を利用することは示されていないが、例えば、ステップS162において、質問文から内容語を抽出して注目語リストに加えれば、質問文を容易に処理することができる。

40 【0262】なお、図50のステップS165、S166において、 $a$ を親話題区間とするすべての話題区間 $b$ を重要文選択の対象とするのではなく、 $b$ として重要な話題区間を1つ選び、それだけを対象として重要文を選択してもよい。この方法は、なるべく短い要約を作成したい場合に有効である。

【0263】例えば、非常に短い要約を作成したい場合には、複数の話題区間から重要文を選択すると、それぞれの話題区間について採録できる量が、読んで理解できる量を下回ってしまうことがある。そのような場合には、要約作成の対象とする話題を絞って要約することで、その話題については理解可能な量の文を採録するこ

とができる。これにより、すべての話題を網羅していても理解するのが困難な要約に比べて、より好ましい要約を作成することができる。

【0264】図52は、このような重要文選択処理のフローチャートである。図52において、ステップS161からステップS164までの処理は、図50と同様である。重要文選択部30は、次に、aを話題区間とする基本窓幅 $w_1$ の話題区間の出現する区間を選択し、それをb0nとする(ステップS165a)。そして、b0nに固有の重要語を、b0nの注目語リスト $kw1n$ に登録する。

【0265】次に、親話題区間aの注目語リスト $kw1a$ の重要語を、 $kw1n$ にマージする(ステップS165b)。そして、b0nを親話題区間とする窓幅 $w_2$ の話題区間の中から、b0nの注目語リスト $kw1a$ に含まれる注目語が最も多く出現する区間を選択し、それを要約対象の話題区間b nとする(ステップS165c)。こうして、1つの親話題区間aから1つの要約対象の話題区間b nを選択した後、図51の処理を行う。

【0266】ここで、1つの親話題区間aからただ1つの要約対象の話題区間b nを選択する代わりに、注目語の出現頻度の大きい順に適当な数の話題区間b nを選択するようにしてもよい。また、要約対象として選択した話題区間b nから十分な量の文を抜粋できない場合には、注目語が次に多く出現する話題区間からも文を選択するようにしてもよい。さらに、ステップS161およびステップS173において、ユーザが入力した質問文の内容語に基づき、特定の親話題区間aのみを処理対象に選ぶことも可能である。

【0267】図53および図54は、図51のステップS167において呼び出される選択処理のフローチャートである。選択処理のサブモジュールは、まず、要約全体の大きさの上限U1と各話題区間の抜粋量の上限U2を、ユーザから受け取る(図53、ステップS181)。通常、U1は、前述の望ましい要約の大きさ $S_a$ より大きく設定され、U2は、前述の望ましい話題の抜粋量 $S_y$ より大きく設定される。これらのパラメータは、 $S_a$ および $S_y$ をもとにして自動的に算出することもできる。

【0268】次に、各話題区間b n毎に、b n内に存在する各文と注目語リスト $kw1n$ 内の注目語とを比較し、注目語の出現数(異なり数と延べ数)を、各文毎に記録する(ステップS182)。そして、U2を越えない長さの未選択の文の中で、注目語の出現数が最大のものを各話題区間b nから1文ずつ選択する(ステップS183)。

【0269】このとき、b n内でそれまでに選択済みの文があれば、それらの長さの和(b nの抜粋量)と新たな選択する文の長さの合計がU2を越えないように、新た

な文を選択する。注目語の出現数としては、異なり数と延べ数のいずれか一方を用いてもよく、両方の合計を用いてもよい。そして、選択された文に選択済であることを示す印を付け、b nの抜粋量に選択された文の長さを加算する。

【0270】次に、文が選択できなかった話題区間b nに選択終了の印を付け(ステップS184)、選択済のすべての文の長さの合計sを求める(ステップS185)。そして、sをU1と比較し(図54、ステップS186)、sがU1以下であれば、処理を終了する。

【0271】s>U1であれば、すべての話題区間b nに選択終了の印を付け(ステップS187)、選択された文の中で注目語の出現数が最小のものを除外して、sとb nの抜粋量をその長さだけ減らす(ステップS188)。そして、再びsをU1と比較し(ステップS189)、まだなおs>U1であれば、sがU1以下になるまでステップS188の処理を繰り返す。

【0272】このような選択処理によれば、最終的に出力される要約文書の大きさは、指定された上限U1以内であることが保証される。上述の要約対象文書の場合、図55、56、および57に示すような要約文書が出力される。ここでは、図面の制約上、1つの要約文書を3つの図に分けて掲載している。この要約文書において、各文の前後に掲げられた記号“...”は、選択されなかった文の存在を示している。

【0273】次に、英語の要約対象文書として、米国出願の明細書の原稿(23, 000語)を用いた例について説明する。ここでは、次のような処理方法およびパラメータを採用した。

(1) 単語認定の方法：ストップワードリストを用いた方法

(2) 結束度計算用の窓の幅：

最大窓幅 $w_0 = 2560$  (語)

基本窓幅 $w_1 = 1280$  (語)

窓幅 $w_2 = 640$  (語)

(3) 話題境界認定の方法：書式ボタンを用いた方法

(4) 重要箇所特定処理における低結束度判定用の感度 $\alpha$ ：

$w_0$  用： $\alpha = -\sigma 0 / 2$  ( $\sigma 0$ は、窓幅 $w_0$ の結束度の標準偏差)

$w_1$  および $w_2$  用： $\alpha = 0$

(5) 重要語抽出の閾値： $h = 6.63490$  (有意水準10%)

(6) 重要文選択における抜粋量の上限値：

U1 = 3, 000 (bytes)

U2 = 600 (bytes)

要約対象文書の全体を掲載することは適当ではないので、参考として、要約対象文書中の見出しの一覧を図58に示す。図58において、()内の表現は、説明のために付加された見出しの省略形であり、要約対象文書に

は含まれていない。

【0274】図59は、入力された要約対象文書の先頭部分を示しており、図60は、その部分に対する単語認定処理の結果を示している。図60において、□で括弧された部分が、認定された単語に対応する。先頭の1文字のみが大文字の単語は、□では、すべて小文字に置き換えられている。

【0275】ここでは、空白および“、”、“.”、“:”、“;”などの区切り記号を手掛かりに単語が切り出され、それらの単語のうち、図61に示すストップワードリストに含まれる単語が取り除かれた。ストップワードリストとは、重要語として抽出したくない冠詞、前置詞などの単語を、あらかじめ定義したリストである。

【0276】また、図62は、図38の統合処理において節境界を求めるために用いた書式ボタンとその境界レベルを示している。ここでは、先頭が大文字のアルファベットで始まっている行を、境界レベル0の節境界とみなし、最初の空白でない文字が“[”である行を、境界レベル1の節境界とみなしている。

【0277】話題境界認定処理においては、話題境界候補区間にこれらの書式ボタンに一致する行が見つかった場合は図38の統合処理を採用し、そうでない場合は図26の統合処理を採用した。その結果、図63に示すような認定結果が得られた。

【0278】図63において、節境界の近くに記された(Bg)、“<1>”などは、図58に示された見出しの省略形を表す。話題境界候補区間データB(1)[p]のうち、“<1>”の節境界P31に対応するB(0)[1]は、書式パターンを用いなければ、B(1)[3]と統合されるべきデータである。ここでは、書式パターンを用いた結果として節境界P31が検出されている。ところが、B(1)およびB(2)の系列に、P31の位置を含むデータが含まれていなかったため、B(1)

[2]およびB(2)[3]のような疑似的な統合対象データが生成されている。

【0279】図64は、重要箇所特定処理の結果を示している。図64において、斜線部分P41およびP42は、窓幅 $w_1$ (1280語)の低結束度区間の除外処理により除外された話題区間を表す。また、横線は、窓幅 $w_0$ の各話題区間における窓幅 $w_1$ の結束度の平均値を表し、矢印は、窓幅 $w_1$ の各話題区間の中心付近において、最大結束度に対応する点を表す。また、ハッチングされた部分P43、P44、およびP45は、窓幅 $w_2$ (640語)の低結束度区間の除外処理により除外された話題区間を表す。

【0280】なお、窓幅 $w_0$ の話題区間に関して低結束度区間の除外処理の感度パラメータ $\alpha$ を前述のように調整したのは、要約対象文書の(Claims)に対応する節の結束度が他の節に比べて異常に高かったためであ

る。このことは、窓幅 $w_0$ の結束度の標準偏差が大きかったことに対応している。実際、窓幅 $w_0$ の結束度の平均値が0.43であるのに対して、その標準偏差は0.11であった。この認定結果に基づき、図65、66、および67に示すような要約文書が生成された。

【0281】以上説明した実施形態においては、日本語および英語の文書を例に挙げて要約処理を説明したが、本発明は、これらの文書以外にも、任意の言語および任意の形式の文書に対して適用され、同様の結果を得ることができる。

【0282】また、要約対象文書は、必ずしもデジタル化された電子文書である必要はなく、例えば、紙媒体などに記載された文書でもよい。この場合、イメージスキャナなどの光電変換装置により文書画像を取り込み、文字認識を行うことで、単語認定可能な文書データを作成することができる。

【0283】

【発明の効果】本発明によれば、数十頁に渡るような長い文書についても、文書サイズの1/2〜1/4程度の大きな話題のまとまりから、段落程度の大きさ(数十語から100語程度)の話題のまとまりまで、任意の大きさの話題のまとまりの階層的構成を、語彙的結束性という文章一般に見られる現象に基づいて認定することができる。

【0284】さらに、それぞれの話題のまとまりから適切な内容を抜粋して、話題の階層的構成に対応する要約を作成することができる。これにより、従来は取り扱いが難しかった、複数の話題に関する文章が混在した複合文書の要約が可能になる。

【0285】また、要約作成の単位とする話題のまとまりの大きさを自動的に決定し、要約対象を重要な話題のまとまりに絞り込むことで、要約として出力すべき大きさに応じて、適切な粒度の話題をバランスよく取り込んだ要約を作成することが可能になる。

【図面の簡単な説明】

【図1】本発明の文書要約装置の原理図である。

【図2】本発明の文書要約装置の構成図である。

【図3】情報処理装置の構成図である。

【図4】記録媒体を示す図である。

【図5】第1の要約対象文書中の見出しを示す図(その1)である。

【図6】第1の要約対象文書中の見出しを示す図(その2)である。

【図7】第1の要約対象文書中の見出しを示す図(その3)である。

【図8】単語認定処理のフローチャートである。

【図9】第1の入力文書を示す図である。

【図10】第1の単語認定結果を示す図である。

【図11】形態素解析処理のフローチャートである。

【図12】日本語の辞書引きの例を示す図である。

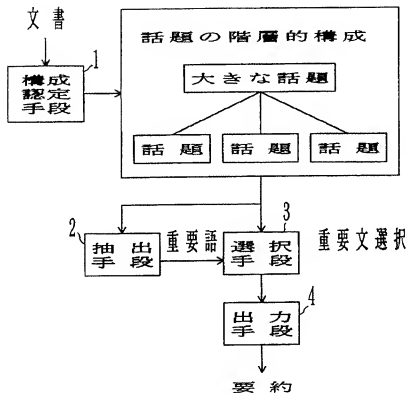


【図 1 3】英語の辞書引きの例を示す図である。  
 【図 1 4】要約粒度決定処理のフローチャートである。  
 【図 1 5】第 1 の結束度分布を示す図である。  
 【図 1 6】第 2 の結束度分布を示す図である。  
 【図 1 7】話題構成認定処理のフローチャートである。  
 【図 1 8】左窓と右窓を示す図である。  
 【図 1 9】窓内の語彙数を示す図である。  
 【図 2 0】結束度の系列を示す図である。  
 【図 2 1】第 3 の結束度分布を示す図である。  
 【図 2 2】移動平均と文書領域の関係を示す図である。  
 【図 2 3】話題境界候補区間認定処理のフローチャートである。  
 【図 2 4】結束力分布を示す図である。  
 【図 2 5】話題境界認定処理のフローチャートである。  
 【図 2 6】第 1 の統合処理のフローチャートである。  
 【図 2 7】統合対象データを示す図である。  
 【図 2 8】疑似データの作成方法を示す図である。  
 【図 2 9】話題構成の第 1 の認定結果を示す図である。  
 【図 3 0】話題構成の第 2 の認定結果を示す図である。  
 【図 3 1】第 1 の話題境界の区間を示す図である。  
 【図 3 2】第 1 の再現率と適合率を示す図である。  
 【図 3 3】第 2 の再現率と適合率を示す図である。  
 【図 3 4】第 2 の話題境界の区間を示す図である。  
 【図 3 5】第 3 の再現率と適合率を示す図である。  
 【図 3 6】第 4 の再現率と適合率を示す図である。  
 【図 3 7】第 1 の書式ボタンと境界レベルを示す図である。  
 【図 3 8】第 2 の統合処理のフローチャートである。  
 【図 3 9】重要箇所特定処理のフローチャート（その 1）である。  
 【図 4 0】重要箇所特定処理のフローチャート（その 2）である。  
 【図 4 1】話題区間除外処理のフローチャートである。  
 【図 4 2】最大結束度計算処理のフローチャートである。  
 【図 4 3】重要箇所の第 1 の特定結果を示す図である。  
 【図 4 4】重要語抽出処理のフローチャート（その 1）である。  
 【図 4 5】重要語抽出処理のフローチャート（その 2）である。  
 【図 4 6】尤度比検定処理のフローチャート（その 1）である。  
 【図 4 7】尤度比検定処理のフローチャート（その 2）である。  
 【図 4 8】話題区間中に含まれている見出しを示す図である。  
 【図 4 9】話題区間から抽出された重要語を示す図である。  
 【図 5 0】重要文選択処理のフローチャート（その 1）である。

【図 5 1】重要文選択処理のフローチャート（その 2）である。  
 【図 5 2】重要文選択処理の他のフローチャートである。  
 【図 5 3】選択処理のフローチャート（その 1）である。  
 【図 5 4】選択処理のフローチャート（その 2）である。  
 【図 5 5】第 1 の要約結果を示す図（その 1）である。  
 【図 5 6】第 1 の要約結果を示す図（その 2）である。  
 【図 5 7】第 1 の要約結果を示す図（その 3）である。  
 【図 5 8】第 2 の要約対象文書中の見出しを示す図である。  
 【図 5 9】第 2 の入力文書を示す図である。  
 【図 6 0】第 2 の単語認定結果を示す図である。  
 【図 6 1】ストップワードを示す図である。  
 【図 6 2】第 2 の書式ボタンと境界レベルを示す図である。  
 【図 6 3】話題構成の第 3 の認定結果を示す図である。  
 【図 6 4】重要箇所の第 2 の特定結果を示す図である。  
 【図 6 5】第 2 の要約結果を示す図（その 1）である。  
 【図 6 6】第 2 の要約結果を示す図（その 2）である。  
 【図 6 7】第 2 の要約結果を示す図（その 3）である。  
 【符号の説明】  
 1 構成認定手段  
 2 抽出手段  
 3 選択手段  
 4 出力手段  
 1 1 要約対象文書  
 1 2 文書要約装置  
 1 3 要約文書  
 2 1 入力部  
 2 2 単語認定部  
 2 3 形態素解析部  
 2 4 単語辞書  
 2 5 要約粒度決定部  
 2 6 話題構成認定部  
 2 7 話題境界候補区間認定部  
 2 8 重要箇所特定部  
 2 9 重要語抽出部  
 3 0 重要文選択部  
 3 1 出力部  
 4 1 出力装置  
 4 2 入力装置  
 4 3 C P U  
 4 4 ネットワーク接続装置  
 4 5 媒体駆動装置  
 4 6 補助記憶装置  
 4 7 メモリ  
 4 8 バス

【図1】

## 本発明の原理図



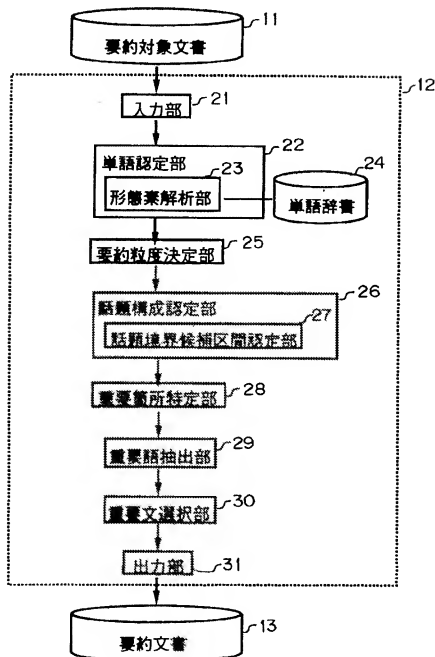
【図12】

日本語の辞書引きの例を示す図

入力文	東京都は大都市だ	
見出し(語幹)	品詞	
東	名詞	
東京都	名詞	
京	名詞	
都	助詞「は」	
は	接頭語	
大	名詞	
都市	助動詞「だ」	
だ		

【図2】

## 文書要約装置の構成図



【図58】

第2の要約対象文書中の  
見出しを示す図

TITLE: SGML Type Document Managing Apparatus and Managing Method

Background of the Invention (Bg)

Field of the Invention

Description of the Related Art

Summary of the Invention

Brief Description of Drawings

Description of Preferred Embodiment (Preferred Embodiment)

[1] Petal editing DTD (< 1 >)

[2] Revision history information (< 2 >)

What is claimed is: (Claim)

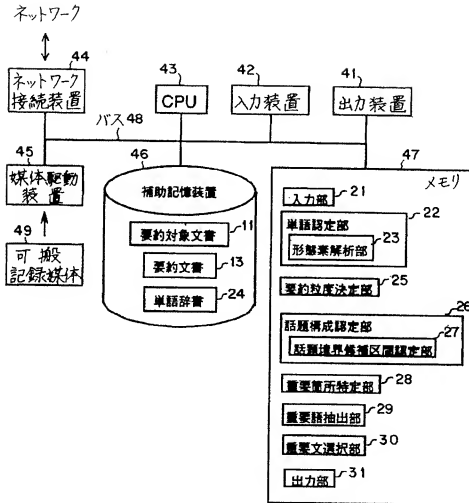
【図62】

第2の書式ボタンと  
境界レベルを示す図

書式ボタン	境界レベル
[A-Z]	0
^_^[	1

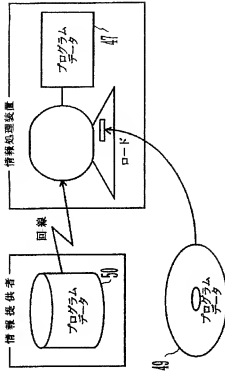
【図3】

## 情報処理装置の構成図



【図4】

記録媒体を示す図



【図31】

第1の話題境界の問題を示す図

話題	認定境界数	話題境界の問題(語)		
		平均	最大値	最低値
2,560	4	3,563	6,145	1,310
1,280	9	1,782	3,010	851
640	17	990	2,185	375
320	38	457	1,060	205
160	74	238	475	40
80	150	111	245	35
40	322	55	185	15

【図9】

第1の入力文書を示す図

インターネットは予想されていた以上の早さで急速に普及している。業務はもちろん特に家庭での利用が急速に広がっている。それにとってもなってインターネットを通じて提供される情報も多種多様化している。そのため、インターネットを利用する上での社会的・技術的なさまざまな要請が顕在化し、それらへの対応が急務となっている。これらの要請の中には、インターネットの健全な運営に関わる認証などの問題とともに、インターネットのサービス内容を高度化するや使い勝手をよくすることに關わる知的情報アクセスの問題が重要な問題として認識されている。

本委員会は、特に、ネットワークアクセスにおける知的情報アクセスの問題について、広範かつ詳細な専門的な立場から調査・研究することを目的としている。その調査・研究では、特に自然言語処理の分野における技術的な側面に兼目し、自然言語処理に関する技術の現状を整理し、かつ将来の発展を見通すことにより、今後において協力してあるいは個別で重点的に取り組むべき課題を明らかにすることを目的としている。

#### 4. ネットワークアクセス技術委員会

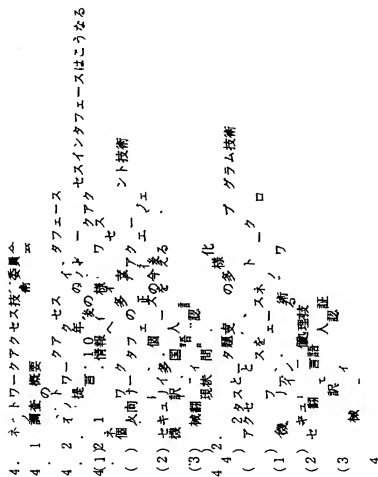
##### 4. 1 調査の概要

インターネットは予想されていた以上の早さで急速に普及している。業務はもちろん特に家庭での利用が急速に広がっている。それにとってもなってインターネットを通じて提供される情報も多種多様化している。そのため、インターネットを利用する上での社会的・技術的なさまざまな要請が顕在化し、それらへの対応が急務となっている。これらの要請の中には、インターネットの健全な運営に関わる認証などの問題とともに、インターネットのサービス内容を高度化するや使い勝手をよくすることに關わる知的情報アクセスの問題が重要な問題として認識されている。

本委員会は、特に、ネットワークアクセスにおける知的情報アクセスの問題について、広範かつ詳細な専門的な立場から調査・研究することを目的としている。その調査・研究では、特に自然言語処理の分野における技術的な側面に兼目し、自然言語処理に関する技術の現状を整理し、かつ将来の発展を見通すことにより、今後において協力してあるいは個別で重点的に取り組むべき課題を明らかにすることを目的としている。

【図5】

第10要約対象文書中の見出しを示す図(その1)



【図6】

第1の要約対象文書中の見出しを示す図(その2)

- 4. 3 ネットワーク上の検索サービス
- 4. 3. 1 検索サービスの調査
  - (1) WWW検索サービスの概要
  - (2) 情報収集/検索方式
  - (3) 情報提示方式
  - (4) 今後の課題
- 4. 3. 2 検索技術の動向
  - (1) キーワード抽出
  - (2) 文書自動分類
  - (3) 要約・抄録技術
  - (4) 分散検索
- 4. 3. 3 電子出版及び電子図書館
  - (1) 電子出版
  - (2) 電子図書館

【図7】

## 第10要約対象文書中の見出しを示す図(その3)

## 4. 4. 検索エンジン

## 4. 4. 1. 日本語の全文検索技術の動向

- (1) 文字列検索アルゴリズム
- (2) インデックス作成法
- (3) 日本語の全文検索技術
- (4) 製品化動向
- (5) 今後の課題

## 4. 4. 2. 有限オートマトンによる自然言語処理技術の動向

- (1) 有限変換器のコンパクト化
- (2) 文字列パタン照合
- (3) 書き換え規則, Two-level モデル
- (4) 形態素解析, 構文解析
- (5) まとめ

## 4. 4. 3 情報フィルタリング技術の動向

- (1) 内容に基づくフィルタリング (content-based filtering)
- (2) 協調フィルタリング (collaborative filtering)
- (3) ユーザモデリング
- (4) まとめ

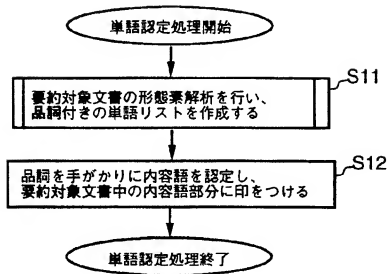
## 4. 4. 4 情報抽出 / 統合技術の動向

- (1) 検索ナビゲーション技術
- (2) 情報統合技術
- (3) 情報の可視化技術



【図8】

## 単語認定処理のフローチャート



【図10】

## 第10単語認定結果を示す図

## 4. 1 【調査/する】の【概要/】

【インターネット/】は【予想/する】されて【い/る】た以上の【早さ/】で【急遽/】に【普及/する】して【い/る】。【業務/】はもちろん特に【業務/】での【利用/する】が【急遽/】に【広/る】って【い/る】。それに【ともな/う】って【インターネット/】を通じて【提供/する】される【情報/】も【多種多様化/】して【い/る】。そのため、【インターネット/】を【利用/する】する上での【社会的/】-【技術的/】なさまざまな【要請/する】が【顕在化/する】し、それらへの【対応/する】が【業務/】と【い/る】。これら【要請/する】の中には、【インターネット/】の【健全/】な【運営/する】に【関わる/る】【認証/する】などの【問題/】とともに、【インターネット/】の【サービス/する/内容/】を【高度化/する】するや【使/う】い【勝手/】をよく【する/】ことに【関わる/る】【知的/】【情報/】【アクセス/する】の【問題/】が【重要/】な【問題/】として【認識/する】されて【い/る】。

本【委員会/】は、特に、【ネットワーク/】【アクセス/する】における【知的/】【情報/】【アクセス/する】の【問題/】について、【広範/】かつ【詳細/】な【専門/】的な【立場/】から【調査/する】・【研究/する】することを【目的/】として【い/る】。その【調査/する】・【研究/する】では、特に、【自然言語/】【処理/する】の【分野/】における【技術的/】な【側面/】に【着目/する】し、【自然言語/】【処理/する】に関する【技術/】の【現状/】を【整理/する】しつつ将来の【発展/する】を【見通/す】すことにより、今後において【協力/する】してあるいは【個別/】で【重点的/】に【取り組/む】むべき【課題/】を【明らかに/】することを【目的/】として【い/る】。

【図32】

第1の再現率と適合率を示す図

範囲	語境界数	一致	再現率	適合率
2,500	2	1	50.0%	25.0%
1,200	3	2	66.7%	22.2%
640	12	5	41.7%	29.4%
320	18	6	33.3%	15.8%
160	33	11	33.3%	14.3%
80	43	15	34.9%	9.4%
40	45	18	40.0%	5.6%

【図33】

第2の再現率と適合率を示す図

範囲	語境界数	一致	再現率	適合率
2,500	3	1	33.3%	33.0%
1,200	12	4	33.3%	44.4%
640	18	6	33.3%	33.3%
320	33	8	24.2%	21.1%
160	43	13	30.2%	17.6%
80	45	15	33.3%	9.4%
40	46	13	41.3%	5.9%

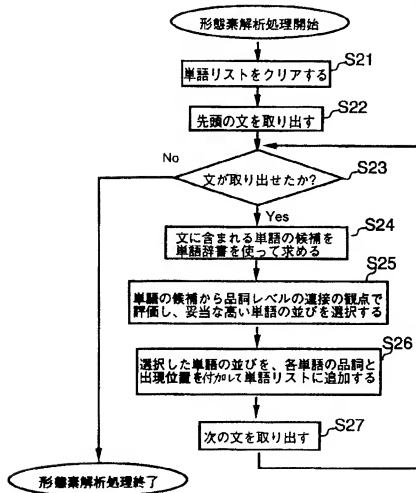
【図34】

第2の語境界の開隔を示す図

範囲	認定境界数	語境界の開隔(語)		
		平均	最大値	最低値
2,500	3	4,454	5,270	2,535
1,200	10	1,620	2,535	805
640	17	950	1,345	605
320	30	575	1,100	170
160	70	251	470	85
80	147	120	230	40
40	308	58	165	10

【図11】

## 形態素解析処理のフローチャート



【図35】

第3の再現率と適合率を示す図

範囲	単語数	一致	再現率	適合率
2,500	2	1	50.0%	33.3%
1,200	3	2	66.7%	20.0%
640	12	5	41.7%	29.4%
320	18	7	38.9%	23.3%
160	33	13	39.4%	18.8%
80	43	18	41.9%	12.2%
40	45	21	46.7%	6.8%

【図13】

英語の辞書引きの例を示す図

入力文	headword	base(root) form	part of speech
Tokyo is the Japanese capital.	Tokyo	Tokyo	proper noun
辞書単語	is	be	be verb (the third person singular present form)
	the	the	definite article
	Japanese	Japanese	proper noun
	Japanese	Japanese	adjective
	capital	capital	noun

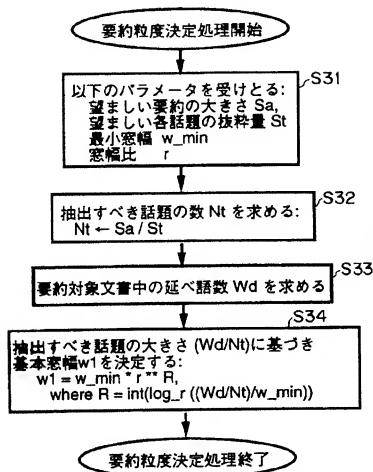
【図19】

窓内の語彙数を示す図

窓幅	40 語
結束度測定位置	40 (語)
左窓(W1)中の語彙数(異なり語数)	29 語
右窓(W2)中の語彙数(異なり語数)	29 語
共通語彙数(異なり語数)	6 語 (い/る, 問題, 調査する, 情報, 関わる, 技術的)

【図 14】

## 要約粒度決定処理のフローチャート



【図 36】

第 4 の再現率と適合率を示す図

窓幅	節境界数	一致	再現率	適合率
2,560	3	1	33.3%	33.3%
1,280	12	4	33.3%	40.0%
640	18	7	38.9%	41.2%
320	38	10	30.3%	33.3%
160	43	14	32.6%	20.0%
80	45	18	40.0%	12.2%
40	46	21	45.7%	6.8%

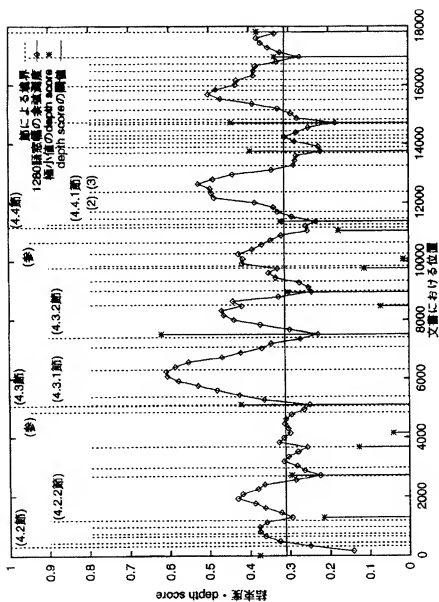
【図 37】

第 1 の  
書式パターンと境界レベルを示す図

書式パターン	境界レベル
$\backslash d + [ \cdot ] + s$	0
$\backslash d + \sim \backslash d + \sim [ \cdot ] + s$	1
$\backslash d + \sim \backslash d + \sim \backslash d + \sim [ \cdot ] + s$	2
$\backslash (d + ) [ \cdot ] + s$	3
$\sim s$	4

【図15】

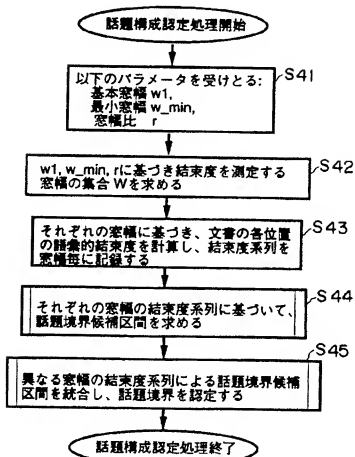
第1の結束度分布を示す図





【図17】

## 話題構成認定処理のフローチャート





【図18】

左窓と右窓を示す図

W1

## 4. 1 【調査する】の【調査】

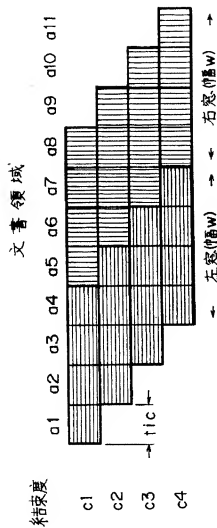
【インターネット】は【予想する】されて【いる】た以上の【早さ】で【急遽】に【普及する】して【い】る。【素部】はもちろん特に【素部】での【利用する】が【急遽】に【広がる】って【い】る。それ【に】と【ともなう】って【インターネット】を通じて【提供する】たれる【情報】も【多種多様化】して【い】る。そのため、【インターネット】を【利用する】する上での【社会的・技術的】な【さまざまな】【要請する】が【顕在化する】し、それらへの【対応する】が【急務】と【なる】って【い】る。これらの【要請する】の中には、【インターネット】の【健全】な【運営する】に【関わる】る【配座する】などの【問題】とともに、【インターネット】の【サービスする】【内容】を【高度化する】するや【使う】い【選手】を【よくする】に【に】に【関わる】る【知的】【情報】【アクセスする】の【問題】が【重要な】な【問題】として【認識する】されて【い】る。

W2

本【委員会】は、特に、【ネットワーク】【アクセスする】における【知的】【情報】【アクセスする】の【問題】について、【広範】かつ【詳細】な【専門】的な【立場】から【調査する】・【研究する】すること【を目的】として【い】る。その【調査する】・【研究する】では、特に、【自然言語】【処理する】の【分野】における【技術的】な【側面】に【着目する】し、【自然言語】【処理する】に【関する】【技術】の【現状】を【整理する】しかつ【将来の】【発展する】を【見通す】することにより、今後において【協力する】して【ある】いは【個別】に【重点的】に【取り組む】むべき【課題】を【明らかに】すること【を目的】として【い】る。

【図20】

結束度の系列を示す図



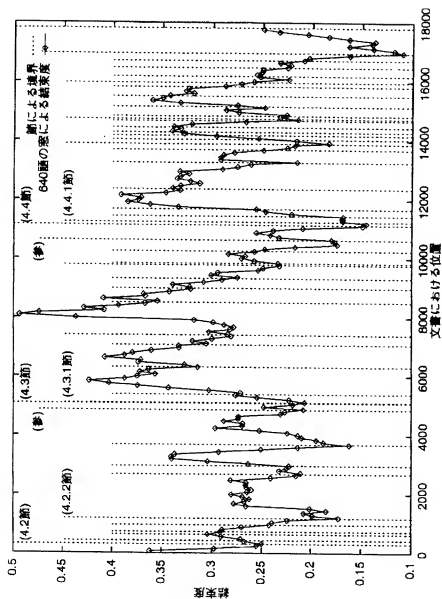
【図22】

移動平均と文書領域の関係を示す図

項目数		文書領域の使用回数										
		a1	a2	a3	a4	a5	a6	a7	a8	a9	a10	a11
4項平均 (c1~c4)	左窓	1	2	3	4	3	2	1	0	0	0	0
	右窓	0	0	0	0	1	2	3	4	3	2	1
3項平均 (c1~c3)	左窓	1	2	3	3	2	1	0	0	0	0	
	右窓	0	0	0	0	1	2	3	3	2	1	
2項平均 (c1, c2)	左窓	1	2	2	2	1	0	0	0	0		
	右窓	0	0	0	0	1	2	2	2	1		

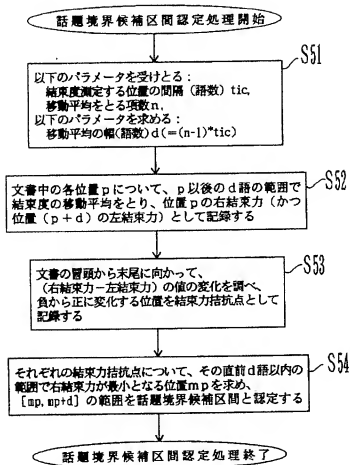
【図21】

第3の結束度分布を示す図



【図 23】

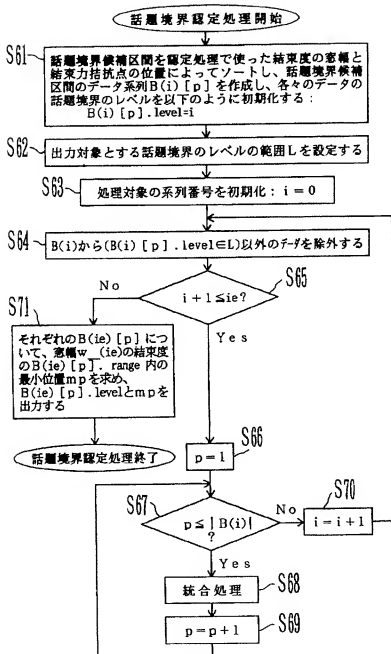
## 話題境界候補区間認定処理のフローチャート





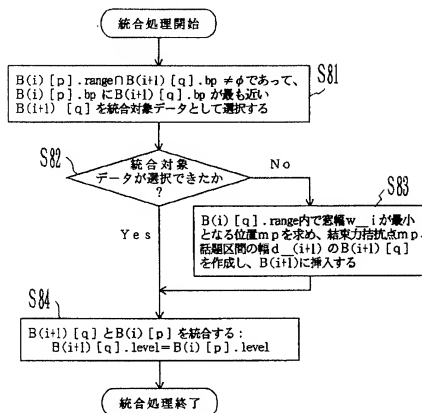
【図25】

## 話題境界認定処理のフローチャート



【図26】

## 第 1 の統合処理のフローチャート



【図48】

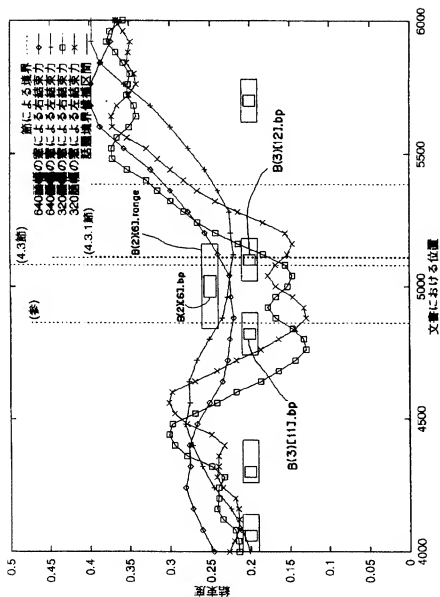
話題区間中に含まれている見出しを示す図

## 話題区間中に含まれている見出し

4. 1 調査の概要
4. 2 ネットワークアクセスのインタフェース
4. 2. 1 提言：10年後のネットワークアクセスインタフェースはこうなる
  - (1) ネットワーク情報への多様なアクセス
  - (2) 個人向けインタフェースを支えるエージェント技術
  - (3) セキュリティ・個人認証の今後
  - (4) 機械翻訳と多国語

【図27】

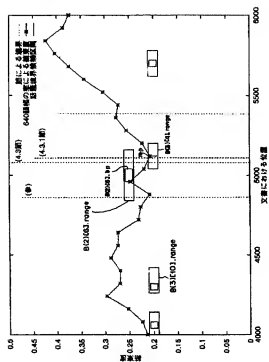
統合対象データを示す図





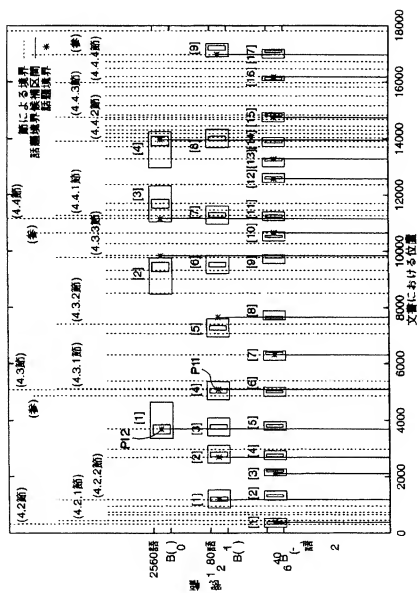
【図 28】

擬似データの作成方法を示す図



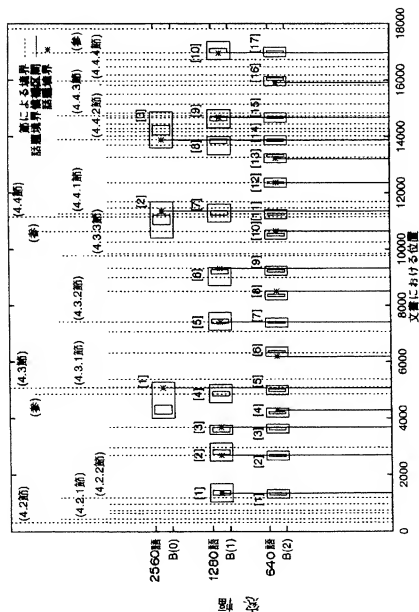
【図29】

話題構成の第1の認定結果を示す図



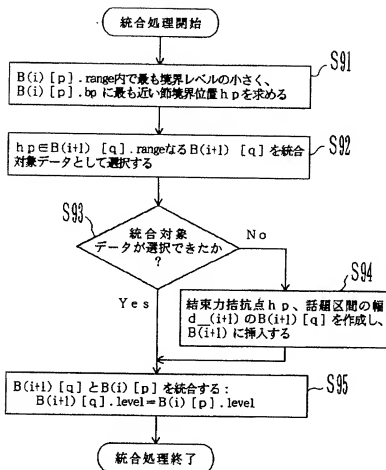
【図30】

話題構成の第2の認定結果を示す図



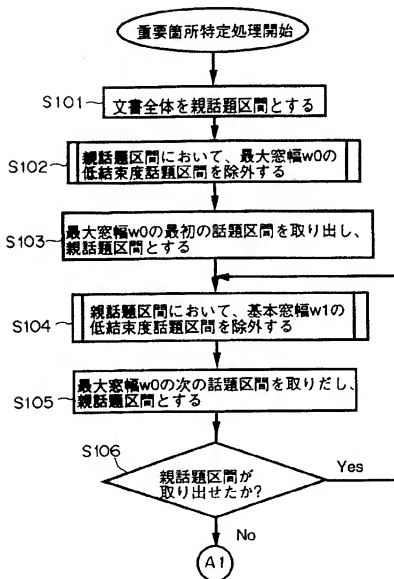
【図38】

## 第 2 の 統 合 処 理 を 示 す 図



【 図 9 】

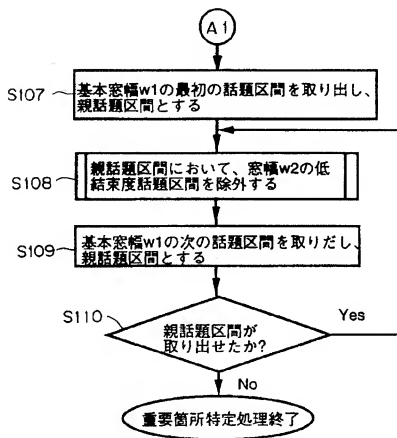
重要箇所特定処理のフローチャート  
( 図 1 )



【図40】

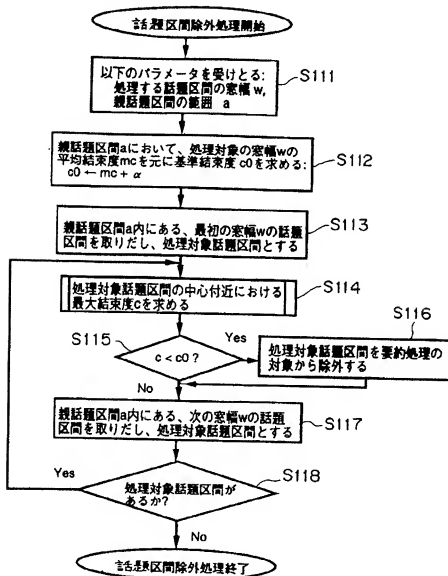
## 重要箇所特定処理のフローチャート

(その2)



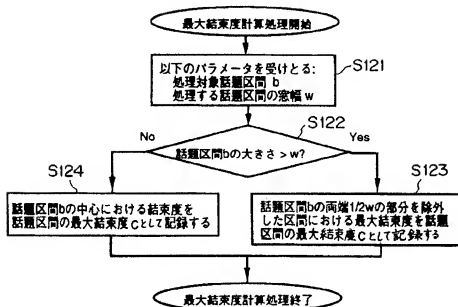
【図 41】

## 話題区間除外処理のフローチャート



【図 4 2】

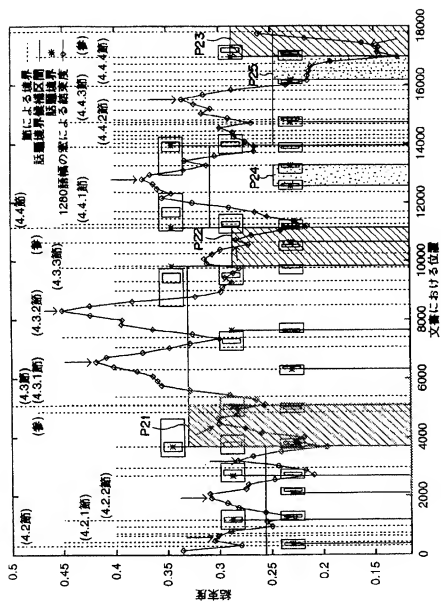
## 最大結束度計算処理のフローチャート





【図43】

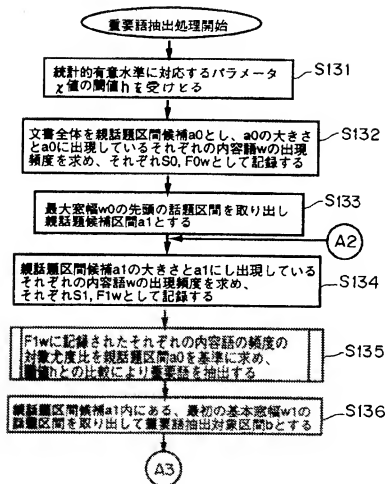
重要箇所の第1の特定結果を示す図



【図44】

## 重要語抽出処理のフローチャート

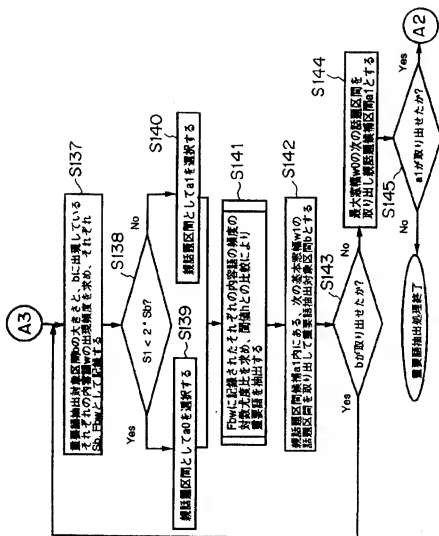
(その1)



【図45】

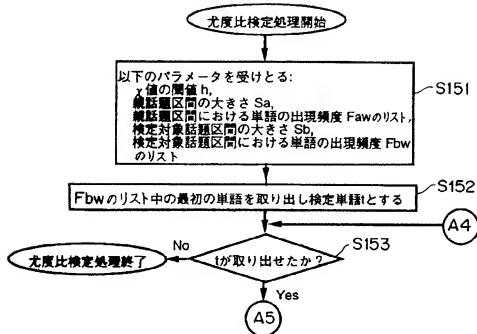
## 重要語抽出処理のフローチャート

(その2)



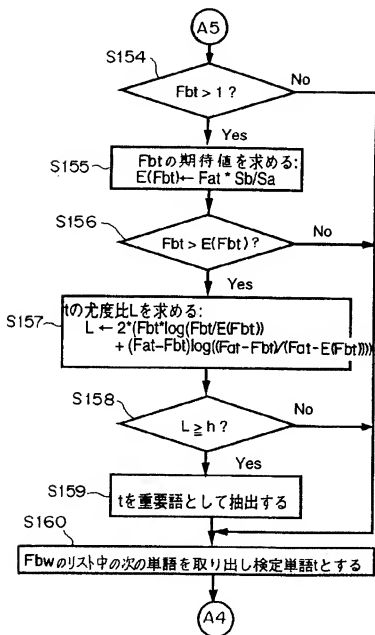
【図 46】

## 尤度比検定処理のフローチャート (その 1)



【図 47】

## 尤度比検定処理のフローチャート(その2)



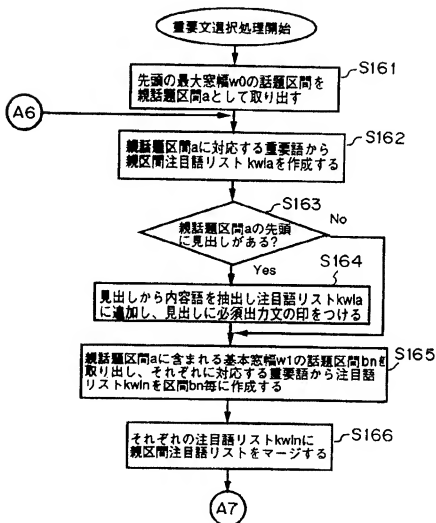
【図49】

話題区間から抽出された重要語を示す図

話題区間から抽出された重要語			
重要語 t	尤度比 L	区間内頻度 Fbt	区間外頻度 Fat
処理 / する	35.6	21	23
技術 /	34.7	38	53
な / る	28.4	34	49
調査 / する	27.3	15	16
アクセス / する	21.2	32	50
インターネット /	21.1	23	32
関わ / る	20.4	9	9
睡眠 / する	18.1	8	8
予想 / する	15.9	7	7
通票 / する	15.9	7	7
自然言語 /	14.7	9	10
インタフェース /	13.7	14	19
ネットワーク /	12.5	21	34
情報 /	11.6	51	106
量 /	11.3	5	5
家庭 /	11.3	5	5
パスワード /	11.3	5	5
重要 /	10.6	7	8
サービス / する	9.75	8	10
管理 / する	9.71	8	10
国語 /	9.08	4	4
母国語 /	9.08	4	4
資源 /	9.08	4	4
発展 / する	9.08	4	4
現状 /	8.66	6	7

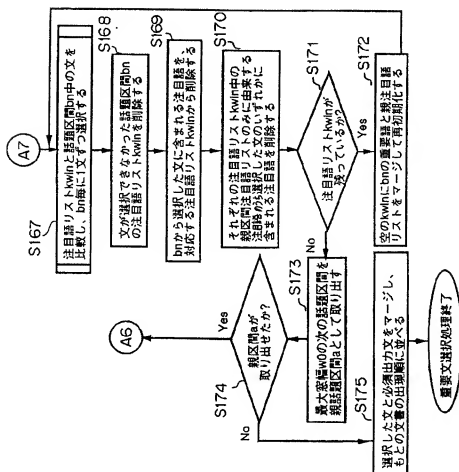
【図50】

## 重要文選択処理のフローチャート（その1）



【図51】

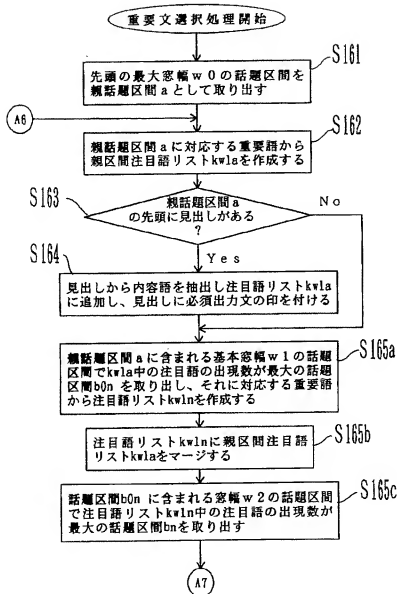
## 重要文選択処理のフローチャート (その2)





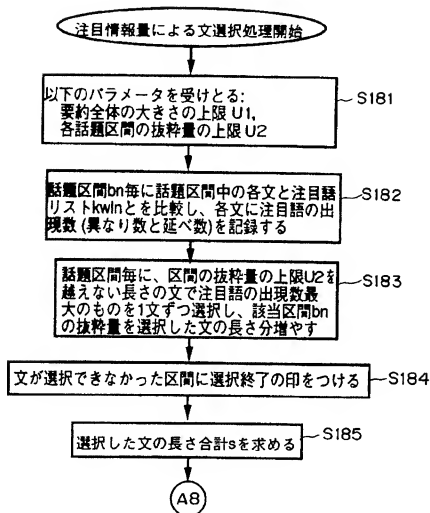
【図52】

## 重要文選択処理の他のフローチャート



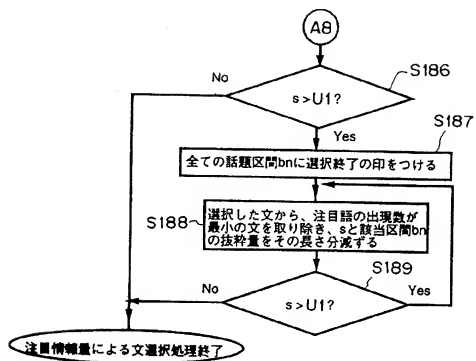
【図53】

## 選択処理のフローチャート (その1)



【図54】

## 選択処理のフローチャート (その2)





【図56】

【図59】

## 第1の要約結果を示す図(その2) 第2の入力文書を示す図

4.3 ネットワーク上の検索サービス  
 ...  
 ...また検索精度を高めるために、高頻度語は検索の対象としない、タイトルや見出しに含まれる語に重みをつける、などの工夫がなされている。  
 ...  
 また、検索サービスが収集したページ数が膨大になるにつれて、ヒット数も膨大になってきたため、すばやく必要な情報を探すために、よりわかりやすい自動抄録作成技術が必要となる。...  
 ...  
 tf・idf方式とは、単語に分割された文章の各単語の重要度を、その単語が文書中に出現する頻度 $tf$ と、その単語を含む文書が文書集合中に出現する頻度の逆数 $idf$ の積によってその単語の重要さを数値化する手法である。  
 ...  
 [河合 92]の研究キーワードのカイ二乗値から各キーワードの分類に対する得点を計算する場合に、シソーラス辞書から得られる抽象的な意味を得点に加える手法である。...  
 ...

## SGML Type Document Managing Apparatus and Managing Method

## Background of the Invention Field of the Invention

The present invention relates to an SGML (Standard Generalized Markup Language) document managing apparatus for allowing users to collaboratively create, edit, and revise a large SGML document sequence, such as a manual.

【図57】

【図60】

第1の辞約結果を示す図 (その3) 第2の単語認定結果を示す図

#### 4.4. 検索エンジン

...  
 [補地、92]では、出現位置を記録する部分文字列検索において、字種に基づく照合単位(漢字1文字、連続するかな2文字、かな文字と非かな文字の接点の2文字など)の設定、文字位置情報の早期配列、低頻度照合単位からの文字位置照合などを組み合わせることににより、専用ハードを用いなくともソフトウェアで日本語の全文検索を高速に行えることを示した。

#### 4.4.2. 有限オートマトンによる自然言語処理技術の動向

...  
 ...しかし、ボタンが有限オートマトンとして与えられた場合は、上記のトライのように状態と変遷(そして、その接尾部分)が一瞥に定まらない為、矢脱遷移が構成できない。...

...  
 ...[Brill, 92]のタガーは、規則ベースでありながら、規則の自動獲得と適用、順序の学習により統計ベースのものと同等レベルの精度が得られ、かつ、よりコンパクトであるという特長があるものの、規則が多くなるにつれ処理速度が低下する。...

SGML[SGML] Type[type] Document[document] Managing [managing] Apparatus[apparatus]  
 and Managing[managing] Method[method]

Background[background] of the Invention[invention] Field[field] of the Invention[invention]  
 The present[present] invention[invention] relates[relates] to an SGML[SGML]  
 (Standard[standard] Generalized[generalized] Markup[markup] Language[language]) docu-  
 ment[document] managing[managing] apparatus[apparatus] for allowing[allowing] users[users]  
 to collaboratively[collaboratively] create [create], edit[edit], and revise[revise] a large[large]  
 SGML[SGML] document[document] sequence[sequence], such as a manual[manual].

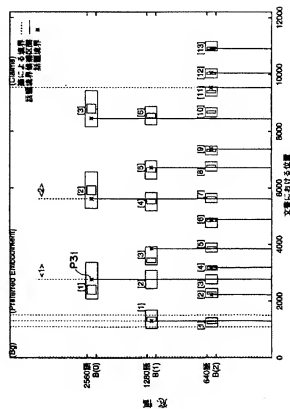
【図61】

ストップワードを示す図

a, after, against, all, along, already, also, although, always, among, an, and, and/or, another, any, anywhere, are, as, at, be, because, been, before, before/after, being, belonging, belongs, below, between, both, but, by, can, cannot, corresponding, do, does, each, either, else, especially, even, every, for, forth, from, further, has, have, he/she, his/her, however, if, in, into, is, it, its, just, later, least, mainly, may, more, moreover, most, much, namely, next, no, not, of, on, once, only, or, other, others, otherwise, out, part, previous, same, should, since, so, some, someone, such, than, that, the, their, them, then, there, thereafter, thereof, these, they, this, those, through, thus, to, too, two, types, under, unless, unlike, until, up, usually, was, well, were, what, when, where, whereby, wherein, whether, which, while, who, whole, whose, why, will, with, with/without, without, yes

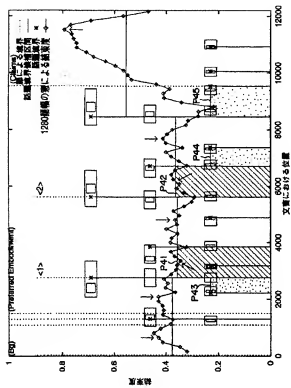
【図63】

話題構成の第3の認定結果を示す図



【図 6 4】

重要箇所第 2 の特定結果を示す図



【図 6 7】

第 2 の要約結果を示す図 (7 の 3)

What is claimed is:

...

15. The document managing apparatus as set forth in claim 14, further comprising:  
editing means for editing the electric document to be edited corresponding to the edited result to said database means when the determined result of said editing consistency examining means is satisfied.

...



【図 6 5】

## 第 2 の要約結果を示す図 (その 1)

## SGML Type Document Managing Apparatus and Managing Method

...

... In addition, SGML documents have been widely used for fields that handle large and long-life documents such as network communications, electronic trading, and databases such as electronic libraries. ...

...

The present invention is mainly intended to provide an SGML type document managing apparatus and an SGML type document managing method that allow collaborative creating and editing works to be effectively performed. ...

...

## Brief Description of Drawings

...

Fig. 2A is a block diagram showing the structure of the client 2 of the information processing system in the case the structure shown in Fig. 1 is accomplished by software.

...

【図66】

## 第2の要約結果を示す図 (その2)

The CPUs 100-1 and 100-2 of the client 2 and the server 3 execute **software** that is stored in the main **storing units** 400-1 and 400-2, the **software** being read from the **auxiliary storing unit** 200-1 or the **input/output units** 500-1 and 500-2, or the **software** being obtained from the **connecting network** through the **network connecting units** 300-1 and 300-2, respectively....

...

At step S71, the SGML document accessing unit 30 adds the fist item of the **repeatable model group** without a head, namely the **removable instance number** and a **generic identifier/model group** with a occurrence indicator, and "(<sup>n</sup>" to the top of the **extracted extended content model**.

...

At step S160, the SGML document editing unit 10 removes the **identifiers** (the document identifier and the **element identifier**) of the **element** stored in the **register** from the **deleted element attribute** (deletions) of the **parent element** of the **paste destination**.

...